



# Red Hat Enterprise Linux for Real Time

## 9

### Understanding RHEL for Real Time

An introduction to fundamental concepts for RHEL for Real Time



# Red Hat Enterprise Linux for Real Time 9 Understanding RHEL for Real Time

---

An introduction to fundamental concepts for RHEL for Real Time

## Legal Notice

Copyright © 2022 Red Hat, Inc.

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at

<http://creativecommons.org/licenses/by-sa/3.0/>

. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, the Red Hat logo, JBoss, OpenShift, Fedora, the Infinity logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux<sup>®</sup> is the registered trademark of Linus Torvalds in the United States and other countries.

Java<sup>®</sup> is a registered trademark of Oracle and/or its affiliates.

XFS<sup>®</sup> is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL<sup>®</sup> is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js<sup>®</sup> is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack<sup>®</sup> Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

## Abstract

This guide helps users understand the fundamental concepts and associated references for optimizing the RHEL for Real Time kernel to maintain low latency, consistent response time, and determinism.

## Table of Contents

|   |           |
|---|-----------|
| <b>MAKING OPEN SOURCE MORE INCLUSIVE</b> .....                                  | <b>4</b>  |
| <b>PROVIDING FEEDBACK ON RED HAT DOCUMENTATION</b> .....                        | <b>5</b>  |
| <b>CHAPTER 1. HARDWARE PLATFORMS FOR RHEL FOR REAL TIME</b> .....               | <b>6</b>  |
| 1.1. PROCESSOR CORES  | 6         |
| 1.2. ADDITIONAL RESOURCES   | 7         |
| <b>CHAPTER 2. MEMORY MANAGEMENT ON RHEL FOR REAL TIME</b> .....                 | <b>8</b>  |
| 2.1. DEMAND PAGING  | 8         |
| 2.2. MAJOR AND MINOR PAGE FAULTS  | 9         |
| 2.3. MLOCK() SYSTEM CALLS   | 10        |
| 2.4. SHARED LIBRARIES   | 11        |
| 2.5. SHARED MEMORY  | 11        |
| <b>CHAPTER 3. HARDWARE INTERRUPTS ON RHEL FOR REAL TIME</b> .....               | <b>13</b> |
| 3.1. LEVEL-SIGNALED INTERRUPTS  | 13        |
| 3.2. MESSAGE-SIGNALED INTERRUPTS  | 13        |
| 3.3. NON-MASKABLE INTERRUPTS  | 14        |
| 3.4. SYSTEM MANAGEMENT INTERRUPTS   | 14        |
| 3.5. ADVANCED PROGRAMMABLE INTERRUPT CONTROLLER                                 | 14        |
| <b>CHAPTER 4. RHEL FOR REAL TIME PROCESSES AND THREADS</b> .....                | <b>16</b> |
| 4.1. PROCESSES  | 16        |
| 4.2. THREADS  | 16        |
| 4.3. ADDITIONAL RESOURCES   | 17        |
| <b>CHAPTER 5. APPLICATION TIMESTAMPING ON RHEL FOR REAL TIME</b> .....          | <b>18</b> |
| 5.1. HARDWARE CLOCKS  | 18        |
| 5.2. POSIX CLOCKS   | 18        |
| 5.3. CLOCK_GETTIME() FUNCTION   | 19        |
| 5.4. ADDITIONAL RESOURCES   | 19        |
| <b>CHAPTER 6. SCHEDULING POLICIES FOR RHEL FOR REAL TIME</b> .....              | <b>20</b> |
| 6.1. SCHEDULER POLICIES   | 20        |
| 6.2. PARAMETERS FOR SCHED_DEADLINE POLICY                                       | 21        |
| <b>CHAPTER 7. AFFINITY IN RHEL FOR REAL TIME</b> .....                          | <b>22</b> |
| 7.1. PROCESSOR AFFINITY   | 22        |
| 7.2. SCHED_DEADLINE AND CPUSETS   | 23        |
| <b>CHAPTER 8. THREAD SYNCHRONIZATION MECHANISMS IN RHEL FOR REAL TIME</b> ..... | <b>24</b> |
| 8.1. MUTEXES  | 24        |
| 8.2. BARRIERS   | 24        |
| 8.3. CONDITION VARIABLES  | 25        |
| 8.4. MUTEX CLASSES  | 25        |
| 8.5. THREAD SYNCHRONIZATION FUNCTIONS   | 26        |
| <b>CHAPTER 9. SOCKET OPTIONS IN RHEL FOR REAL TIME</b> .....                    | <b>28</b> |
| 9.1. TCP_NODELAY SOCKET OPTION  | 28        |
| 9.2. TCP_CORK SOCKET OPTION   | 28        |
| 9.3. EXAMPLE PROGRAMS USING SOCKET OPTIONS                                      | 29        |
| <b>CHAPTER 10. RHEL FOR REAL TIME SCHEDULER</b> .....                           | <b>31</b> |

|   |           |
|---|-----------|
| 10.1. CHRT UTILITY FOR SETTING THE SCHEDULER                | 31        |
| 10.2. PREEMPTIVE SCHEDULING                                 | 31        |
| 10.3. LIBRARY FUNCTIONS FOR SCHEDULER PRIORITY              | 31        |
| <b>CHAPTER 11. SYSTEM CALLS IN RHEL FOR REAL TIME</b> ..... | <b>33</b> |
| 11.1. SCHED_YIELD() FUNCTION                                | 33        |
| 11.2. GETRUSAGE() FUNCTION                                  | 33        |



## MAKING OPEN SOURCE MORE INCLUSIVE

Red Hat is committed to replacing problematic language in our code, documentation, and web properties. We are beginning with these four terms: master, slave, blacklist, and whitelist. Because of the enormity of this endeavor, these changes will be implemented gradually over several upcoming releases. For more details, see [our CTO Chris Wright's message](#).



# PROVIDING FEEDBACK ON RED HAT DOCUMENTATION

We appreciate your feedback on our documentation. Let us know how we can improve it.

## Submitting comments on specific passages

1. View the documentation in the **Multi-page HTML** format and ensure that you see the **Feedback** button in the upper right corner after the page fully loads.
2. Use your cursor to highlight the part of the text that you want to comment on.
3. Click the **Add Feedback** button that appears near the highlighted text.
4. Add your feedback and click **Submit**.

## Submitting feedback through Bugzilla (account required)

1. Log in to the [Bugzilla](#) website.
2. Select the correct version from the **Version** menu.
3. Enter a descriptive title in the **Summary** field.
4. Enter your suggestion for improvement in the **Description** field. Include links to the relevant parts of the documentation.
5. Click **Submit Bug**.

# CHAPTER 1. HARDWARE PLATFORMS FOR RHEL FOR REAL TIME

Configuring the hardware correctly plays a critical role in setting up the real-time environment because hardware impacts the way your system operates. Not all hardware platforms are real-time capable and enable fine tuning. Before performing fine tuning, you must ensure that the potential hardware platform is real-time capable.

Hardware platforms vary based on the vendor. You can test and verify the hardware suitability for real-time with the hardware latency detector (**hwlatdetect**) program. The program controls the latency detector kernel module and helps to detect latencies caused by underlying hardware or firmware behavior.

Any tuning steps required for low latency operation have been completed. Refer to the vendor documentation for instructions to

## Prerequisites

- The RHEL-RT packages are installed.
- Any tuning steps required for low latency operation are complete. Refer to the vendor documentation for instructions to reduce or remove any System Management Interrupts (SMIs) that transition the system into System Management Mode (SMM).



### WARNING

Red Hat recommends not disabling System Management Interrupts (SMIs) completely, as it can result in catastrophic hardware failures.

## 1.1. PROCESSOR CORES

A real-time processor core is a physical Central Processing Unit (CPU) and it executes the machine code. A socket is a connection between the processor and the motherboard of the computer. The socket is the location on the motherboard that the processor is placed into. There are two sets of processors:

- Single core processor that occupies one socket with one available core.
- Quad-core processor that occupies one socket with four available cores.

When designing a real time environment, be aware of the number of available cores, the cache layout among cores, and how the cores are physically connected.

When multiple cores are available, use threads or processes. A program when written without using these constructs, runs on a single processor at a time. A multi-core platform provides advantages through using different cores for different types of operations.

## Caches

Caches have a noticeable impact on overall processing time and determinism. Often, the threads of an application need to synchronize access to a shared resource, such as a data structure.

With the **tuna** command line tool (CLI), you can determine the cache layout and bind interacting threads to a core so that they share the cache. Cache sharing reduces memory faults by ensuring that the mutual exclusion primitive (mutex, condition variables, or similar) and the data structure use the same cache.

## Interconnects

Increasing the number of cores on systems can cause conflicting demands on the interconnects. This makes it necessary to determine the interconnect topology to help detect the conflicts that occur between the cores on real-time systems.

Many hardware vendors now provide a transparent network of interconnects between cores and memory, known as Non-uniform memory access (NUMA) architecture.

NUMA is a system memory design used in multiprocessing, where the memory access time depends on the memory location relative to the processor. When you use NUMA, a processor can access its own local memory faster than non-local memory, such as memory on another processor or memory shared between processors. On NUMA systems, understanding the interconnect topology helps to place threads that communicate frequently on adjacent cores.

The **taskset** and **numactl** utilities determine the CPU topology. **taskset** defines the CPU affinity without NUMA resources such as memory nodes and **numactl** controls the NUMA policy for processes and shared memory.

## 1.2. ADDITIONAL RESOURCES

- [Installing RHEL 9 for Real Time](#)

## CHAPTER 2. MEMORY MANAGEMENT ON RHEL FOR REAL TIME

Real-time systems use a virtual memory system, where an address referenced by a user-space application translates into a physical address. The translation occurs through a combination of page tables and address translation hardware in the underlying computing system. An advantage of having the translation mechanism in between a program and the actual memory is that the operating system can swap pages when required or on request of the CPU.

In real-time, to swap pages from the secondary storage to the primary memory, the previously used page table entry is marked as invalid. As a result, even under normal memory pressure, the operating system can retrieve pages from one application and give them to another. This can cause unpredictable system behaviors.

Memory allocation implementations include demand paging mechanism and memory lock (**mlock()**) system calls.



### NOTE

Sharing data information on CPUs in different cache and NUMA domains might cause traffic problems and bottlenecks.

When writing a multithreaded application, consider the machine topology before designing the data decomposition. Topology is the memory hierarchy, and includes CPU caches and the Non-Uniform Memory Access (NUMA) node.

### 2.1. DEMAND PAGING

Demand paging is similar to a paging system with page swapping. The system loads pages that are stored in the secondary memory when required or on CPU demand. All memory addresses generated by a program pass through an address translation mechanism in the processor. The addresses then convert from a process-specific virtual address to a physical memory address. This is referred to as virtual memory. The two main components in the translation mechanism are page tables and translation lookaside buffers (TLBs)

#### Page tables

Page tables are multi-level tables in physical memory that contain mappings for virtual to physical memory. These mappings are readable by the virtual memory translation hardware in the processor.

A page table entry with an assigned physical address, is referred to as the resident working set. When the operating system needs to free memory for other processes, it can swap pages from the resident working set. When swapping pages, any reference to a virtual address within that page creates a page fault and causes page reallocation.

When the system is extremely low on physical memory, the swap process starts to thrash, which constantly steals pages from processes, and never allows a process to complete. You can monitor the virtual memory statistics by looking for the **pgfault** value in the **/proc/vmstat** file.

#### Translation lookaside buffers

Translation Lookaside Buffers (TLBs) are hardware caches of virtual memory translations. Any processor core with a TLB checks the TLB in parallel with initiating a memory read of a page table entry. If the TLB entry for a virtual address is valid, the memory read is aborted and the value in the TLB is used for the address translation.

A TLB operates on the principle of locality of reference. This means that if code stays in one region of memory for a significant period of time (such as loops or call-related functions) then the TLB references avoid the main memory for address translations. This can significantly speed up processing times.

When writing deterministic and fast code, use functions that maintain locality of reference. This can mean using loops rather than recursion. If recursion are not avoidable, place the recursion call at the end of the function. This is called tail-recursion, which makes the code work in a relatively small region of memory and avoids calling table translations from the main memory.

## 2.2. MAJOR AND MINOR PAGE FAULTS

RHEL for Real Time allocates memory by breaking physical memory into chunks called pages and then maps them to the virtual memory. Faults in real-time occur when a process needs a specific page that is not mapped or is no longer available in the memory. Hence faults essentially mean unavailability of pages when required by a CPU. When a process encounters a page fault, all threads freeze until the kernel handles this fault. There are several ways to address this problem, but the best solution can be to adjust the source code to avoid page faults.

### Minor page faults

Minor page faults in real-time occur when a process attempts to access a portion of memory before it has been initialized. In such scenarios, the system performs operations to fill the memory maps or other management structures. The severity of a minor page fault can depend on the system load and other factors, but they are usually short and have a negligible impact.

### Major page faults

Real-time major faults occur when the system has to either synchronize memory buffers with the disk, swap memory pages belonging to other processes, or undertake any other input output (I/O) activity to free memory. This occurs when the processor references a virtual memory address that does not have a physical page allocated to it. The reference to an empty page causes the processor to execute a fault, and instructs the kernel code to allocate a page which increases latency dramatically.

In real-time, when an application shows a performance drop, it is beneficial to check for the process information related to page faults in the `/proc/` directory. For a specific process identifier (PID), using the `cat` command, you can view the `/proc/PID/stat` file for following relevant entries:

- Field 2: the executable file name.
- Field 10: the number of minor page faults.
- Field 12: the number of major page faults.

The following example demonstrates viewing page faults with the `cat` command and a `pipe` function to return only the second, tenth, and twelfth lines of the `/proc/PID/stat` file:

```
# cat /proc/3366/stat | cut -d\ -f2,10,12
(bash) 5389 0
```

In the example output, the process with PID 3366 is bash and it has 5389 minor page faults and no major page faults.

### Additional resources

- Linux System Programming by Robert Love

## 2.3. MLOCK() SYSTEM CALLS

The memory lock (**mlock()**) system calls enable calling processes to lock or unlock a specified range of the address space and prevents Linux from paging the locked memory to swap space. After you allocate the physical page to the page table entry, references to that page are relatively fast. The memory lock system calls fall into **mlock()** and **munlock()** categories.

The **mlock()** and **munlock()** system calls lock and unlock a specified range of process address pages. When successful, the pages in the specified range remain resident in the memory until the **munlock()** system call unlocks the pages.

The **mlock()** and **munlock()** system calls take following parameters:

- **addr**: specifies the start of an address range.
- **len**: specifies the length of the address space in bytes.

When successful, **mlock()** and **munlock()** system calls return 0. In case of an error, they return -1 and set a **errno** to indicate the error.

The **mlockall()** and **munlockall()** system calls locks or unlocks all of the program space.



### NOTE

The **mlock()** system call does not ensure that the program will not have a page I/O. It ensures that the data stays in the memory but cannot ensure it stays on the same page. Other functions such as **move\_pages** and memory compactors can move data around regardless of the use of **mlock()**.

Memory locks are made on a page basis and do not stack. If two dynamically allocated memory segments share the same page locked twice by **mlock()** or **mlockall()**, they unlock by using a single **munlock()** or **munlockall()** system call. As such, it is important to be aware of the pages that the application unlocks to avoid double-locking or single-unlocking problems.

The following are two most common workarounds to mitigate double-lock or single-unlock problems:

- Tracking the allocated and locked memory areas and creating a wrapper function that verifies the number of page allocations before unlocking a page. This is the resource counting principle used in device drivers.
- Making memory allocations based on page size and alignment to avoid double-lock on a page.

### Additional resources

- **capabilities(7)** man page
- **mlock(2)** man page
- **mlock(3)** man page
- **mlockall(2)** man page
- **mmap(2)** man page
- **move\_pages(2)** man page

- **posix\_memalign(3)** man page
- **posix\_memalign(3p)** man page

## 2.4. SHARED LIBRARIES

The RHEL for Real Time shared libraries are called dynamic shared objects (DSO) and are a collection of pre-compiled code blocks called functions. These functions are reusable in multiple programs and they load at run-time or compile time.

Linux supports the following two library classes:

- Dynamic or shared libraries: exists as separate files outside of the executable file. These files load into the memory and get mapped at run-time.
- Static libraries: are files linked to a program statically at compile time.

The **ld.so** dynamic linker loads the shared libraries required by a program and then executes the code. The DSO functions load the libraries in the memory once and multiple processes can then reference the objects by mapping into the address space of processes. You can configure the dynamic libraries to load at compile time using the **LD\_BIND\_NOW** variable.

Evaluating symbols before program initialization can improve performance because evaluating at application run-time can cause latency if the memory pages are located on an external disk.

### Additional resources

- **ld.so(8)** man page

## 2.5. SHARED MEMORY

In RHEL for Real Time, shared memory is a memory space shared between multiple processes. Using program threads, all threads created in one process context can share the same address space. This makes all data structures accessible to threads. With POSIX shared memory calls, you can configure processes to share a part of the address space.

You can use the following supported POSIX shared memory calls:

- **shm\_open()**: creates and opens a new or opens an existing POSIX shared memory object.
- **shm\_unlink()**: unlinks POSIX shared memory objects.
- **mmap()**: creates a new mapping in the virtual address space of the calling process.



### NOTE

The mechanism for sharing a memory region between two processes using System V IPC **shmsem()** set of calls is deprecated and is no longer supported on RHEL for Real Time.

### Additional resources

- **shm\_open(3)** man page
- **shm\_overview(7)** man page

- **mmap(2)** man page



## CHAPTER 3. HARDWARE INTERRUPTS ON RHEL FOR REAL TIME

Real-time systems receive many interrupts over the course of its operation, including a semi-regular "timer" interrupt that periodically performs maintenance and system scheduling decisions. The systems may also receive special kinds of interrupts, such as Nonmaskable Interrupt (NMI) and System Management Interrupt (SMI). Hardware interrupts are used by devices to indicate a change in the physical state of the system that requires attention. For example, a hard disk signaling that it has read a series of data blocks, or when a network device has processed a buffer containing network packets.

When an interrupt occurs in real-time, the system stops active programs are stopped and executes an interrupt handler is executed.

In real-time, hardware interrupts are referenced by an interrupt number. These numbers are mapped back to the piece of hardware that created the interrupt. This enables the system to monitor which device created the interrupt and when it occurred. When an interrupt occurs in real-time, the system stops active programs and executes an interrupt handler. The handler preempts other running programs and system activities. This can slow the entire system and create latencies.

RHEL for Real Time modifies the way interrupts are handled in order to improve performance and decrease latency. Using the **cat /proc/interrupts** command you can print an output to view the types of hardware interrupts that took place, the number of interrupts received, the target CPU for the interrupt, and the device generating the interrupt.

### 3.1. LEVEL-SIGNALED INTERRUPTS

In real-time, level-signaled interrupts use a dedicated interrupt line that delivers voltage transitions. The device controller raises an interrupt by asserting a signal on the interrupt request line. The interrupt line sends one of two voltages to represent a binary 1 or binary 0.

When the interrupt signal is sent by the line, it remains in that state until the CPU resets it. The CPU performs a state save, captures the interrupt, and dispatches the interrupt handler. The interrupt handler determines the cause of the interrupt, clears the interrupt by performing necessary services, and restores the state of the device. The Level-signaled interrupts are more reliable and supports multiple devices, though they are complex to implement.

### 3.2. MESSAGE-SIGNALED INTERRUPTS

In real-time, many systems use message-signaled interrupts (MSI), which send the signal as a dedicated message on a packet or message-based electrical bus. A common example of this type of bus is the Peripheral Component Interconnect Express (PCI Express or PCIe). These devices transmit a message type, which the PCIe host controller interprets as an interrupt message. The host controller then sends the message on to the CPU.

In real-time, depending on the hardware, a PCIe system does one of the following:

- Sends the signal using a dedicated interrupt line between the PCIe host controller and the CPU.
- Sends the message over the CPU HyperTransport bus.

In real-time, PCIe systems can also operate in legacy mode, where legacy interrupt lines are implemented in order to support older operating systems or on boot Linux kernels with the option **pci=noms**i on the kernel command line.

### 3.3. NON-MASKABLE INTERRUPTS

In real-time, non-maskable interrupts are hardware interrupts that standard interrupt masking techniques in the system cannot ignore. The NMIs have higher priority than the maskable interrupts. The NMIs occur to signal attention for non recoverable hardware errors.

In real-time, NMIs are also used by some systems as a hardware monitor. When a processor receives NMI, it handles the NMI immediately by calling the NMI handler pointed to by the interrupt vector. If certain conditions are met, such as an interrupt not being triggered after a specified length of time, the NMI handler signals a warning and provides debugging information about the problem. This helps to identify and prevent system lockups.

In real-time, maskable interrupts are hardware interrupts that can be ignored by setting a bit in an interrupt mask register's bit-mask. CPUs can temporarily ignore maskable interrupts during critical processing.

### 3.4. SYSTEM MANAGEMENT INTERRUPTS

In real-time, system management interrupts (SMIs) offer extended functionality, such as legacy hardware device emulation and can also be used for system management tasks. SMIs are similar to non maskable interrupts (NMIs) in that they use a special electrical signalling line and are generally not maskable. When an SMI occurs, the CPU enters the System Management Mode (SMM). In this mode, a special low-level handler executes to handle the SMIs. The SMM is typically provided directly from the system management firmware, often the BIOS or the EFI.

The real-time SMIs are most often used to provide legacy hardware emulation. A common example is to imitate a diskette drive. When there is no diskette drive attached, the operating system attempts to access the diskette and triggers a SMI. In this scenario, a handler provides the operating system with an emulated device instead. The operating system then treats the emulation as a legacy device.

In real-time, SMIs can adversely affect the system because they take place without the direct involvement of the operating system. A poorly written SMI handling routine may consume many milliseconds of CPU time, and the operating system might not be able to preempt the handler. This can create periodic high latencies in an otherwise well-tuned and highly responsive system. As a vendor may use SMI handlers to manage CPU temperature and fan control, it may not be possible to disable them. In such situations, you must notify the vendor of problems that occur when using these interrupts.

In real-time, you can isolate SMIs using the **hwlatdetect** utility. It is available in the **rt-tests** package. This utility measures the time period during which the CPU is used by an SMI handling routine.

### 3.5. ADVANCED PROGRAMMABLE INTERRUPT CONTROLLER

The advanced programmable interrupt controller (APIC) developed by Intel Corporation, provides the ability to:

- Handle large amounts of interrupts to route each to a specific set of CPUs.
- Support inter-CPU communication and remove the need for multiple devices to share a single interrupt line.

The real-time APIC represents a series of devices and technologies that work together to generate, route, and handle a large number of hardware interrupts in a scalable and manageable way. It uses a combination of a local APIC built into each system CPU and a number of Input/Output APICs that are connected directly to hardware devices.

In real-time, when a hardware device generates an interrupt, the connected I/O APIC detects and routes the interrupt across the system APIC bus to a specific CPU. The operating system knows the IO-APIC connects to the devices, and interrupt line within that device. The Advanced Configuration and Power Interface Differentiated System Description Table (ACPI DSDT) includes information about the specific wiring of the host system motherboard and peripheral components and a device provides information on the available interrupt sources. Together, these two sets of data provide information about the overall interrupt hierarchy.

RHEL for Real Time supports Complex APIC-based interrupt management strategies with the system APICs connected in hierarchies and delivering interrupts to CPUs in a load-balanced fashion rather than targeting a specific CPU or set of CPUs.

## CHAPTER 4. RHEL FOR REAL TIME PROCESSES AND THREADS

The RHEL for Real Time key factors in operating systems are minimal interrupt latency and minimal thread switching latency. Although all programs use threads and processes, RHEL for Real Time handles them in a different way compared to the standard Red Hat Enterprise Linux.

In real-time, using parallelism helps achieve greater efficiency in task execution and latency. Parallelism is when multiple tasks or several sub-tasks run at the same time using the multi-core infrastructure of CPU.

### 4.1. PROCESSES

A real-time process, in simplest terms, is a program in execution. The term process refers to an independent address space, potentially containing multiple threads. When the concept of more than one process running inside one address space was developed, Linux turned to a process structure that shares an address space with another process. This works well, as long as the process data structure is small.

A UNIX®-style process construct contains:

- Address mappings for virtual memory.
- An execution context (PC, stack, registers).
- State and accounting information.

In real-time, each process starts with a single thread, often called the parent thread. You can create additional threads from parent threads using the **fork()** system calls. **fork()** creates a new child process which is identical to the parent process except for the new process identifier. The child process runs independent of the creating process. The parent and child processes can be executed simultaneously. The difference between the **fork()** and **exec()** system calls is that, **fork()** starts a new process which is the copy of the parent process and **exec()** replaces the current process image with the new one.

In real-time, the **fork()** system call, when successful, returns the process identifier of the child process and the parent process returns a non-zero value. On error, it returns an error number.

### 4.2. THREADS

In real-time, multiple threads can exist within a process. All threads of a process share its virtual address space and system resources. A thread is a schedulable entity that contains:

- A program counter (PC).
- A register context.
- A stack pointer.

In real-time, following are potential mechanisms to create parallelism:

- Using the **fork()** and **exec()** function calls to create new processes. The **fork()** call creates an exact duplicate of a process from which it is called and has a unique process identifier.
- Using the Posix threads (**pthread**s) API to create new threads within an already running process.

You must evaluate the component interaction level before forking real-time threads. Creating a new address space and running it as a new process is beneficial when the components are independent of one another or with less interaction. When components are required to share data or communicate frequently, running the threads within one address space is more efficient.

In real-time, the **fork()** system call, when successful, returns a zero value. On error, it returns an error number.

### 4.3. ADDITIONAL RESOURCES

- **fork(2)** man page
- **exec(2)** man page

## CHAPTER 5. APPLICATION TIMESTAMPING ON RHEL FOR REAL TIME

Applications that perform frequent **timestamps** are affected by the CPU cost of reading the clock. The high cost and amount of time used to read the clock can have a negative impact on an application's performance.

You can reduce the cost of reading the clock by selecting a hardware clock that has a reading mechanism, faster than that of the default clock.

In RHEL for Real Time, a further performance gain can be acquired by using POSIX clocks with the **clock\_gettime()** function to produce clock readings with the lowest possible CPU cost.

These benefits are more evident on systems which use hardware clocks with high reading costs.

### 5.1. HARDWARE CLOCKS

Multiple instances of clock sources found in multiprocessor systems, such as non-uniform memory access (NUMA) and Symmetric multiprocessing (SMP), interact among themselves and the way they react to system events, such as CPU frequency scaling or entering energy economy modes, determine whether they are suitable clock sources for the real-time kernel.

The preferred clock source is the Time Stamp Counter (TSC). If the TSC is not available, the High Precision Event Timer (HPET) is the second best option. However, not all systems have HPET clocks, and some HPET clocks can be unreliable.

In the absence of TSC and HPET, other options include the ACPI Power Management Timer (ACPI\_PM), the Programmable Interval Timer (PIT), and the Real Time Clock (RTC). The last two options are either costly to read or have a low resolution (time granularity), therefore they are sub-optimal for use with the real-time kernel.

### 5.2. POSIX CLOCKS

POSIX is a standard for implementing and representing time sources. You can assign a POSIX clock to an application without affecting other applications in the system. This is in contrast to hardware clocks which are selected by the kernel and implemented across the system.

The function used to read a given POSIX clock is **clock\_gettime()**, which is defined at `<time.h>`. The kernel counterpart to **clock\_gettime()** is a system call. When a user process calls **clock\_gettime()**:

1. The corresponding C library (**glibc**) calls the **sys\_clock\_gettime()** system call.
2. **sys\_clock\_gettime()** performs the requested operation.
3. **sys\_clock\_gettime()** returns the result to the user program.

However, the context switch from the user application to the kernel has a CPU cost. Even though this cost is very low, if the operation is repeated thousands of times, the accumulated cost can have an impact on the overall performance of the application. To avoid context switching to the kernel, thus making it faster to read the clock, support for the **CLOCK\_MONOTONIC\_COARSE** and **CLOCK\_REALTIME\_COARSE** POSIX clocks was added, in the form of a virtual dynamic shared object (VDSO) library function.

Time readings performed by **clock\_gettime()**, using one of the **\_COARSE** clock variants, do not require kernel intervention and are executed entirely in user space. This yields a significant performance gain.

Time readings for **\_COARSE** clocks have a millisecond (ms) resolution, meaning that time intervals smaller than 1 ms are not recorded. The **\_COARSE** variants of the POSIX clocks are suitable for any application that can accommodate millisecond clock resolution.



#### NOTE

To compare the cost and resolution of reading POSIX clocks with and without the **\_COARSE** prefix, see the [RHEL for Real Time Reference guide](#).

### 5.3. CLOCK\_GETTIME() FUNCTION

The following code shows an example of code using the **clock\_gettime()** function with the **CLOCK\_MONOTONIC\_COARSE** POSIX clock:

```
#include <time.h>
main()
{
    int rc;
    long i;
    struct timespec ts;

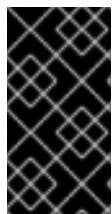
    for(i=0; i<100000000; i++) {
        rc = clock_gettime(CLOCK_MONOTONIC_COARSE, &ts);
    }
}
```

You can improve upon the example above by adding checks to verify the return code of **clock\_gettime()**, to verify the value of the **rc** variable, or to ensure the content of the **ts** structure is to be trusted.



#### NOTE

The **clock\_gettime()** man page provides more information on writing more reliable applications.



#### IMPORTANT

Programs using the **clock\_gettime()** function must be linked with the **rt** library by adding **-lrt** to the **gcc** command line.

```
$ gcc clock_timing.c -o clock_timing -lrt
```

### 5.4. ADDITIONAL RESOURCES

- **clock\_gettime()** man page
- Linux System Programming by Robert Love
- Understanding The Linux Kernel by Daniel P. Bovet and Marco Cesati

## CHAPTER 6. SCHEDULING POLICIES FOR RHEL FOR REAL TIME

In real-time, the scheduler is the kernel component that determines the runnable thread to execute. Each thread has an associated scheduling policy and a static scheduling priority (**sched\_priority**). The scheduling being preemptive, the currently running thread stops when a thread with a higher static priority gets ready to execute. The current thread then returns to the **waitlist** for its static priority.

All Linux threads have one of the following scheduling policies:

- **SCHED\_OTHER** or `^SCHED_NORMAL`: the default policy
- **SCHED\_BATCH**: similar to **SCHED\_OTHER**, but with a throughput orientation
- **SCHED\_IDLE**: a lower priority policy than **SCHED\_OTHER**.
- **SCHED\_FIFO**: a first in and first out real-time policy.
- **SCHED\_RR**: a round-robin real-time policy.
- **SCHED\_DEADLINE**: a scheduler policy to prioritize tasks according to its job deadline, the earliest absolute deadline runs first.

### 6.1. SCHEDULER POLICIES

In real-time, threads have higher priority than the normal threads. The policies have scheduling priority values that range from the minimum value of 1 to the maximum value of 99.

The following policies are critical to real-time:

- **SCHED\_OTHER** or **SCHED\_NORMAL**:  
This is the default scheduling policy for Linux threads. It has a dynamic priority that is changed by the system based on the characteristics of the thread. **SCHED\_OTHER** threads have nice values between 20 (highest priority) and 19 (lowest priority). By default, **SCHED\_OTHER** threads have a nice value of 0.
- **SCHED\_FIFO** policy  
Threads with **SCHED\_FIFO**, run ahead of **SCHED\_OTHER** tasks. Instead of using nice values, **SCHED\_FIFO** uses a fixed priority between 1 (lowest) and 99 (highest). A **SCHED\_FIFO** thread with a priority of 1 always schedules ahead of any **SCHED\_OTHER** thread.
- **SCHED\_RR** policy:  
The **SCHED\_RR** policy is similar to the **SCHED\_FIFO** policy. The threads of equal priority are scheduled in a round-robin fashion. **SCHED\_FIFO** and **SCHED\_RR** threads run until one of the following events occurs:
  - The thread goes to sleep or waits for an event.
  - A higher-priority real-time thread gets ready to run.  
Unless one of the above events occur, the threads run indefinitely on the specified processor, while the lower-priority threads remain in the queue waiting to run. This can cause the system service threads to be resident and prevent being swapped out and fail the filesystem data flushing.



- **SCHED\_DEADLINE**: **SCHED\_DEADLINE** policy specifies the timing requirements. It schedules each task according to the task's deadline, the task with the earliest deadline runs first. The kernel requires **runtime<=deadline<=period** to be true. The relation between the required options is **runtime<=deadline<=period**.

## 6.2. PARAMETERS FOR SCHED\_DEADLINE POLICY

Each **SCHED\_DEADLINE** task is characterized by **period**, **runtime**, and **deadline** parameters. The values for these parameters are integers of nanoseconds. The following table lists these parameters.

Table 6.1. SCHED\_DEADLINE parameters

| Parameter       | Description  |
|-----------------|--|
| <b>period</b>   | <p><b>period</b> is the activation pattern of a real-time task.</p> <p>For example, if a video processing task has 60 frames per second to process, a new frame will be queued for service every 16 milliseconds, so the period is 16 milliseconds.</p>  |
| <b>runtime</b>  | <p><b>runtime</b> is the amount of CPU execution time allotted to the task to produce output. In real-time, the maximum execution time, also known as "Worst Case Execution Time" (WCET) is the <b>runtime</b>.</p> <p>For example, a video processing tool may take, in the worst case, five milliseconds to process an image. Hence its run time is five milliseconds.</p> |
| <b>deadline</b> | <p><b>deadline</b> is the maximum time for the output to be produced.</p> <p>For example, if the task needs to deliver the processed frame within ten milliseconds, the <b>deadline</b> is ten milliseconds.</p>   |

## CHAPTER 7. AFFINITY IN RHEL FOR REAL TIME

In real-time, every thread and interrupt source in the system has a processor affinity property. The operating system scheduler uses this information to determine which threads and interrupts to run on which CPU.

The Affinity in real-time, is represented as a bitmask, where each bit in the mask represents a CPU core. If the bit is set to 1, then the thread or interrupt may run on that core; if 0 then the thread or interrupt is excluded from running on the core. The default value for an affinity bitmask is all ones, meaning the thread or interrupt may run on any core in the system.

By default, processes can run on any CPU. However, processes can be instructed to run on a predetermined selection of CPUs, by changing the affinity of the process. Child processes inherit the CPU affinities of their parents.

Some of the more typical affinity setups include:

- Reserve one CPU core for all system processes and allow the application to run on the remainder of the cores.
- Allow a thread application and a given kernel thread (such as the network **softirq** or a driver thread) on the same CPU.
- Pair producer and consumer threads on each CPU.



### NOTE

The affinity settings must be designed in conjunction with the program for good expected behavior.

## 7.1. PROCESSOR AFFINITY

In real-time, the processes by default, can run on any CPU. However, you can configure the processes to run on a predetermined selection of CPUs, by changing the affinity of the process. Child processes inherit the CPU affinities of their parents.

The real-time practice for tuning affinities on a system is to determine the number of cores required to run the application and then isolating those cores. This can be achieved with the Tuna tool, or with shell scripts to modify the bitmask value.

The **taskset** command can be used to change the affinity of a process and modifying the **/proc/** filesystem entry changes the affinity of an interrupt. Using the **taskset** command with the **-p** or **--pid** option and the process identifier (PID) of the process, checks the affinity of a process.

The **-c** or **--cpu-list** option displays the numerical list of cores, instead of as a bitmask. The affinity can be set by specifying the number of the CPU to bind a specific process. For example, for a process that previously used either CPU 0 or CPU 1, you can change the affinity so that it can only run on CPU 1. In addition to the **taskset** command, you can also set the processor affinity can using the **sched\_setaffinity()** system call.

### Additional resources

- **taskset(1)** man page
- **sched\_setaffinity(2)** man page

## 7.2. SCHED\_DEADLINE AND CPUSSETS

The kernel's deadline scheduling class (**SCHED\_DEADLINE**) implements early deadline first scheduler (EDF) for sporadic tasks with a constrained deadline. It prioritizes the tasks according to the job deadline: earliest absolute deadline first. In addition to the EDF scheduler, the deadline scheduler also implements the constant bandwidth server (CBS). The CBS algorithm is a resource reservation protocol.

The CBS guarantees that each task receives its run time (**Q**) at every period (**T**). At the start of every activation of a task, the CBS replenishes the task's run time. As the job runs, it consumes its **runtime** and if the task runs out of its **runtime**, the task is throttled and de-scheduled. The throttling mechanism prevents a single task from running more than its runtime and helps to avoid the performance problems of other jobs.

In real-time, to avoid the overloading the system with **deadline** tasks, the **deadline** scheduler implements an acceptance test, which is run every time a task is configured to run with the **deadline scheduler**. The acceptance test guarantees that **SCHED\_DEADLINE** tasks does not use more CPU time than the specified on the **kernel.sched\_rt\_runtime\_us/kernel.sched\_rt\_period\_us** files, which is 950 ms over 1s, by default.

## CHAPTER 8. THREAD SYNCHRONIZATION MECHANISMS IN RHEL FOR REAL TIME

In real-time, when two or more threads need access to a shared resource at the same time, the threads coordinate using the thread synchronization mechanism. Thread synchronization ensures that only one thread uses the shared resource at a time. The three thread synchronization mechanisms used on Linux: Mutexes, Barriers, and Condition variables (**condvars**).

### 8.1. MUTEXES

Mutex derives from the terms mutual exclusion. The mutual exclusion object synchronizes access to a resource. It is a mechanism that ensures only one thread can acquire a mutex at a time.

The **mutex** algorithm creates a serial access to each section of code, so that only one thread executes the code at any one time. Mutexes are created using an attribute object known as the **mutex** attribute object. It is an abstract object, which contains several attributes that depends on the POSIX options you choose to implement. The attribute object is defined with the **pthread\_mutex\_t** variable. The object stores the attributes defined for the mutex. The **pthread\_mutex\_init(&my\_mutex, &my\_mutex\_attr)**, **pthread\_mutexattr\_setrobust()** and **pthread\_mutexattr\_getrobust()** functions return 0, when successful. On error, they return the error number.

In real-time, you can either retain the attribute object to initialize more mutexes of the same type or you can clean up (destroy) the attribute object. The mutex is not affected in either case. Mutexes include the standard and advanced type of mutexes.

#### Standard mutexes

The real-time standard mutexes are private, non-recursive, non-robust, and non-priority inheritance capable mutexes. Initializing a **pthread\_mutex\_t** using **pthread\_mutex\_init(&my\_mutex, &my\_mutex\_attr)** creates a standard mutex. When using the standard mutex type, your application may not benefit from the advantages provided by the  **pthreads**  API and the RHEL for Real Time kernel.

#### Advanced mutexes

Mutexes defined with additional capabilities are called advanced mutexes. Advanced capabilities include priority inheritance, robust behavior of a mutex, and shared and private mutexes. For example, for robust mutexes, initializing the **pthread\_mutexattr\_setrobust()** function, sets the robust attribute. Similarly, using the attribute **PTHREAD\_PROCESS\_SHARED**, allows any thread to operate on the mutex, provided the thread has access to its allocated memory. The attribute **PTHREAD\_PROCESS\_PRIVATE** sets a private mutex.

A non-robust mutex does not release automatically and stays locked until you manually release it.

#### Additional resources

- **futex(7)** man page
- **pthread\_mutex\_destroy(P)** man page

### 8.2. BARRIERS

Barriers operate in a very different way when compared to other thread synchronization methods. The barriers define a point in the code where all active threads stop until all threads and processes reach this barrier. Barriers are used in situations when a running application needs to ensure that all threads have completed specific tasks before execution can continue.

The barrier mutex in real-time, take following two variables:

- The first variable records the **stop** and **pass** state of the barrier.
- The second variable records the total number of threads that enter the barrier.

The barrier sets the state to **pass** only when the specified number of threads reach the defined barrier. When the barrier state is set to **pass**, the threads and processes proceed further. The **pthread\_barrier\_init()** function allocates the required resources to use the defined barrier and initializes it with the attributes referenced by the **attr** attribute object.

The **pthread\_barrier\_init()** and **pthread\_barrier\_destroy()** functions return the zero value, when successful. On error, they return an error number.

### 8.3. CONDITION VARIABLES

In real-time, condition variables (**condvar**) is a POSIX thread construct that waits for a particular condition to be achieved before proceeding. In general, the signaled condition relates to the state of data that the thread shares with another thread. For example, a **condvar** can be used to signal a data entry into a processing queue and a thread waiting to process that data from the queue. Using the **pthread\_cond\_init()** function, you can initialize a condition variable.

The **pthread\_cond\_init()**, **pthread\_cond\_wait()**, and **pthread\_cond\_signal()** functions return the zero value, when successful. On error, it returns the error number.

### 8.4. MUTEX CLASSES

The following table lists the mutex classes to consider when writing or porting an application.

**Table 8.1. Mutex options**

| Advanced mutexes               | Description   |
|--------------------------------|---|
| Shared mutexes                 | Defines shared access for multiple threads to acquire a mutex at a given time. Shared mutexes can create a high additional overhead. The attribute is <b>PTHREAD_PROCESS_SHARED</b> .   |
| Private mutexes                | Ensures that only the threads created within the same process can access the mutex. The attribute is <b>PTHREAD_PROCESS_PRIVATE</b> .   |
| Real-time priority inheritance | Sets the priority level of the lower priority task higher above a current higher priority task. When the task completes, it releases the resource and the task drops back to its original priority permitting the higher priority task to execute. The attribute is <b>PTHREAD_PRIO_INHERIT</b> . |

| Advanced mutexes | Description  |
|------------------|--|
| Robust mutexes   | Sets the robust mutexes to release automatically when the owning thread terminates. The value substring <b>NP</b> in the string <b>PTHREAD_MUTEX_ROBUST_NP</b> , indicates that robust mutexes are non-POSIX or not portable |

#### Additional resources

- **futex(7)** man page

## 8.5. THREAD SYNCHRONIZATION FUNCTIONS

The following tables lists the functions applicable for thread synchronization mechanisms

Table 8.2. Functions

| Function   | Description  |
|--|--|
| <b>pthread_mutexattr_init(&amp;my_mutex_attr)</b>    | Initiates a mutex with attributes specified by <b>attr</b> . If <b>attr</b> is NULL, it applies the default mutex attributes.  |
| <b>pthread_mutexattr_destroy(&amp;my_mutex_attr)</b> | Destroys the specified mutex object. You can re-initialize with <b>pthread_mutex_init()</b> .  |
| <b>pthread_mutexattr_setrobust(t)</b>                | Specifies the <b>PTHREAD_MUTEX_ROBUST</b> attribute of a mutex. The <b>PTHREAD_MUTEX_ROBUST</b> attribute defines a thread to terminate without unlocking the mutex. A future call to own this mutex succeeds automatically and returns the value <b>EOWNERDEAD</b> to indicate that the previous mutex owner no longer exists.  |
| <b>pthread_mutexattr_getrobust(t)</b>                | Queries the <b>PTHREAD_MUTEX_ROBUST</b> attribute of a mutex.  |
| <b>pthread_barrier_init()</b>                        | Allocates the required resources to use and initialize the barrier with attribute object <b>attr</b> . If <b>attr</b> is NULL, it applies the default values.  |
| <b>pthread_cond_init()</b>                           | Initializes a condition variable. The <b>cond</b> argument defines the object to initiate with the attributes in the condition variable attribute object <b>attr</b> . If <b>attr</b> is NULL, it applies the default values.  |
| <b>pthread_cond_wait()</b>                           | Blocks a thread execution until it receives a signal from another thread. In addition, a call to this function also releases the associated lock on mutex before blocking. The argument <b>cond</b> defines the <b>pthread_cond_t</b> object for a thread to block on. The <b>mutex</b> argument specifies the mutex to unblock. |

| Function                     | Description   |
|------------------------------|---|
| <b>pthread_cond_signal()</b> | Unblocks at least one of the threads that are blocked on a specified condition variable. The argument <b>cond</b> specifies using the <b>pthread_cond_t</b> object to unblock the thread. |

## CHAPTER 9. SOCKET OPTIONS IN RHEL FOR REAL TIME

The real-time socket is a two way data transfer mechanism between two processes on same systems such as the UNIX domain and loopback devices or on different systems such as network sockets.

Transmission Control Protocol (TCP) is the most common transport protocol and is often used to achieve consistent low latency for a service that requires constant communication or to cork the sockets in a low priority restricted environment.

With new applications, hardware features, and kernel architecture optimizations, TCP has to introduce new approaches to handle the changes effectively. The new approaches can cause unstable program behaviors. Because the program behavior changes as the underlying operating system components change, they must be handled with care.

One example of such behavior in TCP is the delay in sending small buffers. This allows sending them as one network packet. Buffering small writes to TCP and sending them all at once generally works well, but it can also create latencies. For real-time applications, the **TCP\_NODELAY** socket option disables the delay and sends small writes as soon as they are ready.

The relevant socket options for data transfer are **TCP\_NODELAY** and **TCP\_CORK**.

### 9.1. TCP\_NODELAY SOCKET OPTION

The **TCP\_NODELAY** socket option disables Nagle's algorithm. Configuring **TCP\_NODELAY** with the **setsockopt** sockets API function sends multiple small buffer writes as individual packets as soon as they are ready.

Sending multiple logically related buffers as a single packet by building a contiguous packet before sending, achieves better latency and performance. Alternatively, if the memory buffers are logically related but not contiguous, you can create an I/O vector and pass it to the kernel using **writenv** on a socket with **TCP\_NODELAY** enabled.

The following example illustrates enabling **TCP\_NODELAY** through the **setsockopt** sockets API.

```
int one = 1;
setsockopt(descriptor, SOL_TCP, TCP_NODELAY, &one, sizeof(one));
```



#### NOTE

To use **TCP\_NODELAY** effectively, avoid small, logically related buffer writes. With **TCP\_NODELAY**, small writes make TCP send multiple buffers as individual packets, which may result in poor overall performance.

#### Additional resources

- **sendfile(2)** man page

### 9.2. TCP\_CORK SOCKET OPTION

The **TCP\_CORK** option collects all data packets in a socket and prevents from transmitting them until the buffer fills to a specified limit. This enables applications to build a packet in the kernel space and send data when **TCP\_CORK** is disabled. **TCP\_CORK** is set on a socket file descriptor using the **setsockopt()** function. When developing programs, if you must send bulk data from a file, consider using **TCP\_CORK** with the **sendfile()** function.



When a logical packet is built in the kernel by various components, enable **TCP\_CORK** by configuring it to a value of 1 using the **setsockopt** sockets API. This is known as "corking the socket". **TCP\_CORK** can cause bugs if the cork is not removed at an appropriate time.

The following example illustrates enabling **TCP\_CORK** through the **setsockopt** sockets API.

```
int one = 1;
setsockopt(descriptor, SOL_TCP, TCP_CORK, &one, sizeof(one));
```

In some environments, if the kernel is not able to identify when to remove the cork, you can manually remove it as follows:

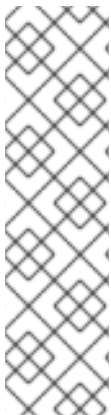
```
int zero = 0;
setsockopt(descriptor, SOL_TCP, TCP_CORK, &zero, sizeof(zero));
```

### Additional resources

- **sendfile(2)** man page

## 9.3. EXAMPLE PROGRAMS USING SOCKET OPTIONS

The **TCP\_NODELAY** and **TCP\_CORK** socket options significantly influence the behavior of a network connection. **TCP\_NODELAY** disables the Nagle's algorithm on applications that benefit by sending data packets as soon as they are ready. With **TCP\_CORK**, you can transfer multiple data packets simultaneously, with no delays between them.



### NOTE

To enable the socket options, for example **TCP\_NODELAY**, build it with the following code and then set appropriate options.

```
gcc tcp_nodelay_client.c -o tcp_nodelay_client -lrt
```

When you run the **tcp\_nodelay\_server** and **tcp\_nodelay\_client** programs without any arguments, the client uses the default socket options. For more information on **tcp\_nodelay\_server** and **tcp\_nodelay\_client** programs, see the [TCP changes result in latency performance when small buffers are used](#) article.

The example programs provide information on the performance impact these socket options can have on your applications.

### Performance impact on a client

You can send small buffer writes to a client without using the **TCP\_NODELAY** and **TCP\_CORK** socket options. When run without any arguments, the client uses the default socket options.

- To initiate data transfer, define the server TCP port and the number of packets it must process. For example, 10,000 packets in this test.

```
$. /tcp_nodelay_server 5001 10000
```

The code sends 15 packets, each of two bytes, and waits for a response from the server. It adopts the default TCP behavior here

## Performance impact on a loopback interface

To enable the socket option, build it using `gcc tcp_nodelay_client.c -o tcp_nodelay_client -lrt` and then set the appropriate options.

Following examples use a loopback interface to demonstrate three variations:

- To send buffer writes immediately, set the **no\_delay** option on a socket configured with **TCP\_NODELAY**.

```
$ ./tcp_nodelay_client localhost 5001 10000 no_delay
```

```
10000 packets of 30 bytes sent in 1649.771240 ms: 181.843399 bytes/ms using
TCP_NODELAY
```

TCP sends the buffers right away, disabling the algorithm that combines the small packets. This improves performance but can cause a flurry of small packets to be sent for each logical packet.

- To collect multiple data packets and send them with one system call, configure the **TCP\_CORK** socket option.

```
$ ./tcp_nodelay_client localhost 5001 10000 cork
```

```
10000 packets of 30 bytes sent in 850.796448 ms: 352.610779 bytes/ms using TCP_CORK
```

Using the cork technique significantly reduces the time required to send data packets as it combines full logical packets in its buffers and sends fewer overall network packets. You must ensure to remove the **cork** at the appropriate time.

When developing programs, if you must send bulk data from a file, consider using **TCP\_CORK** with the **sendfile()** option.

- To measure performance without using socket options.

```
$ ./tcp_nodelay_client localhost 5001 10000
```

```
10000 packets of 30 bytes sent in 400129.781250 ms: 0.749757 bytes/ms
```

This is the baseline measure when TCP combines buffer writes and waits to check for more data than can optimally fit in the network packet.

## Additional resources

- **sendfile(2)** man page

## CHAPTER 10. RHEL FOR REAL TIME SCHEDULER

RHEL for Real Time uses the command line utilities help you to configure and monitor process configurations.

### 10.1. CHRT UTILITY FOR SETTING THE SCHEDULER

The **chrt** utility checks and adjusts scheduler policies and priorities. It can start new processes with the desired properties, or change the current properties of a running process.

The **chrt** utility takes the either **--pid** or the **-p** option to specify the process ID (PID).

The **chrt** utility takes the following policy options:

- **-f** or **--fifo**: sets the schedule to **SCHED\_FIFO**.
- **-o** or **--other**: sets the schedule to **SCHED\_OTHER**.
- **-r** or **--rr**: sets schedule to **SCHED\_RR**.
- **-d** or **--deadline**: sets schedule to **SCHED\_DEADLINE**.

The following example shows the attributes for a specified process.

```
# chrt -p 468
pid 468's current scheduling policy: SCHED_FIFO
pid 468's current scheduling priority: 85
```

### 10.2. PREEMPTIVE SCHEDULING

The real-time preemption is the mechanism to temporarily interrupt an executing task, with the intention of resuming it at a later time. It occurs when a higher priority process interrupts the CPU usage. Preemption can have a particularly negative impact on performance, and constant preemption can lead to a state known as thrashing. This problem occurs when processes are constantly preempted and no process ever gets to run completely. Changing the priority of a task can help reduce involuntary preemption.

You can check for voluntary and involuntary preemption occurring on a single process by viewing the contents of the **/proc/PID/status** file, where PID is the process identifier.

The following example shows the preemption status of a process with PID 1000.

```
# grep voluntary /proc/1000/status
voluntary_ctxt_switches: 194529
nonvoluntary_ctxt_switches: 195338
```

### 10.3. LIBRARY FUNCTIONS FOR SCHEDULER PRIORITY

The real-time processes use a different set of library calls to control policy and priority. The functions require the inclusion of the **sched.h** header file. The symbols **SCHED\_OTHER**, **SCHED\_RR** and **SCHED\_FIFO** must also be defined in the **sched.h** header file.

The following table lists the functions that set the policy and priority for the real-time scheduler.

Table 10.1. Library functions for real-time scheduler

| Functions                       | Description   |
|---------------------------------|---|
| <b>sched_getscheduler()</b>     | Retrieves the scheduler policy for a specific process identifier (PID)  |
| <b>sched_setscheduler()</b>     | Sets the scheduler policy and other parameters. This function requires three parameters:<br><b>sched_setscheduler(pid_t pid, int policy, const struct sched_param *sp);</b> |
| <b>sched_getparam()</b>         | Retrieves the scheduling parameters of a scheduling policy.   |
| <b>sched_setparam()</b>         | Sets the parameters associated with a scheduling policy that has been already set and can be verified using the <b>sched_getparam()</b> function.                           |
| <b>sched_get_priority_max()</b> | Returns the maximum valid priority associated with the scheduling policy.   |
| <b>sched_get_priority_min()</b> | Returns the minimum valid priority associated with the scheduling policy .  |
| <b>sched_rr_get_interval()</b>  | Displays the allocated <b>timeslice</b> for each process.   |

## CHAPTER 11. SYSTEM CALLS IN RHEL FOR REAL TIME

The real-time system call is a function used by application programs to communicate with the kernel. It is a mechanism for programs to order resources from the kernel.

### 11.1. SCHED\_YIELD() FUNCTION

The **sched\_yield()** function is designed for a processor to select a process other than the running one. This type of request is prone to failure when issued from within a poorly-written application.

When the **sched\_yield()** function is used within processes with real-time priorities, it can display unexpected behavior. The process that calls **sched\_yield()** moves to the tail of the queue of processes running at same priority. When there are no other processes running at the same priority, the process that called **sched\_yield()** continues to run. If the priority of that process is high, it can potentially create a busy loop, rendering the machine unusable.

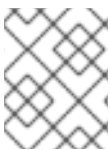
In general, do not use **sched\_yield()** on real-time processes.

### 11.2. GETRUSAGE() FUNCTION

The **getrusage()** function retrieves important information from a specified process or its threads. It reports on information such as:

- The number of voluntary and involuntary context switches.
- Major and minor page faults.
- Amount of memory in use.

**getrusage()** enables you to query an application to provide information relevant to both performance tuning and debugging activities. **getrusage()** retrieves information that would otherwise need to be cataloged from several different files in the **/proc/** directory and would be hard to synchronize with specific actions or events on the application.



#### NOTE

Not all the fields contained in the structure filled with **getrusage()** results are set by the kernel. Some of them are kept for compatibility reasons only.

#### Additional resources

- **getrusage(2)** man page