

Real-time Linux Scheduling Latency

after S.Rostedt, D.B. Oliveira, D. Casini, R.Oliveira,
T.Cucinotta

March 2021

- 1 Real-Time
 - Definition ?
 - Applications ?
 - Approaches with Linux ?
 - PREEMPT_RT patches
 - rt-tests suite and cyclictst
- 2 Model of IRQs, NMIs and thread synchronization mechanisms
- 3 Model the latency bound
- 4 Conclusion

1) Real-Time

Real-Time (RT) :

- is about timing behavior not performance
- deterministic/predictable scheduling
- ~ opposite of “batch work” on a server (think of a build server where performance is about long-run rather than reactivity)
- ⇒ timing response guarantees/bounding = safety bound
- ⇒ predict the worst case

With Linux ? :

- understanding timing behavior of linux
- priority-base scheduling : high priority needs to be able to preempt low priority
- ⇒ faster in worst case scenarios but slower in the average scenarios (otherwise this would be the default kernel)

Applications :

- robotics
- stock exchange
- music studio recording (no “glitches” with jack low-latency audio server)
- death or life devices ?
 - ⇒ NO ! Linux kernel too complex for formal methods ?

Two approaches with the Linux kernel :

- dual-kernel approach : RTAI, RTLinux, **Xenomai Cobalt**, Xenomai Mercure
 - ⇒ Linux becomes a task alongside high-priority RT tasks
 - ⇒ lots of work : support new architecture, implement specific tools/libraries (libc)
 - ⇒ bad scaling, (wo)man power problem
 - in-kernel approach
 - ⇒ maximize preemptible code sections = allow scheduling almost everywhere
 - ⇒ take advantage of all the tools/optimizations/industrial support. . .
 - ⇒ Linux too big for < 1\$ chip ?
 - ⇒ try to keep up-to-date with mainline kernel
 - ⇒ 2021 : try to merge the patches in mainline kernel
- ⇒ both Linux-approaches are ~15/20 years old

Standard Linux preemption model configuration :

PREEMPT_NONE, default = **PREEMPT_VOLUNTARY**,
PREEMPT :

```
linuxconfig qemu.cfs.debug.busybox.config - Linux/x86 5.11.0 Kernel Configuration
> General setup

Preemption Model
Use the arrow keys to navigate this window or press the
hotkey of the item you wish to select followed by the <SPACE
BAR>. Press <?> for additional information about this

( X ) No Forced Preemption (Server)
( ) Voluntary Kernel Preemption (Desktop)
( ) Preemptible Kernel (Low-Latency Desktop)

< .elect > < Help >
```

With patches (new **PREEMPT_RT**) :

```
linuxconfig qemu.cfs.debug.busybox.rt.config - Linux/x86 5.10.17 Kernel Configuration
> General setup

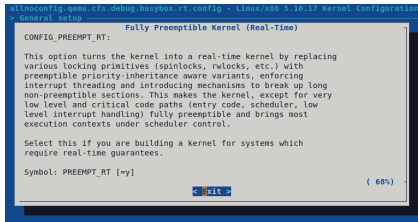
Preemption Model
Use the arrow keys to navigate this window or press the
hotkey of the item you wish to select followed by the <SPACE
BAR>. Press <?> for additional information about this

( ) No Forced Preemption (Server)
( ) Voluntary Kernel Preemption (Desktop)
( ) Preemptible Kernel (Low-Latency Desktop)
( X ) Fully Preemptible Kernel (Real-Time)

< .elect > < Help >
```

Debian binary : linux-image-amd64 ⇔ linux-image-rt-amd64

Make preemption enabled almost everywhere :



```
allnoconfig.gnu.cfs.debug.busybox.rt.config - Linux/x86 5.10.17 Kernel Configuration
> General setup
  CONFIG_PREEMPT_RT: Fully Preemptible Kernel (Real-Time)

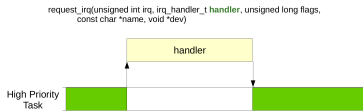
  This option turns the kernel into a real-time kernel by replacing
  various locking primitives (spinlocks, rwlocks, etc.) with
  preemptible priority-inheritance aware variants, enforcing
  interrupt threading and introducing mechanisms to break up long
  non-preemptible sections. This makes the kernel, except for very
  low level and critical code paths (entry code, scheduler, low
  level interrupt handling) fully preemptible and brings most
  execution contexts under scheduler control.

  Select this if you are building a kernel for systems which
  require real-time guarantees.

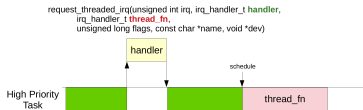
  Symbol: PREEMPT_RT [-y] ( 68%)
```

- spinlocks \rightsquigarrow raw_spinlock + sleeping spinlocks (= rt mutex)
 - threaded interrupt handler
 - priority inheritance to avoid priority inversion = when a high priority is blocked because of some task of lower priority
 - others hacks already merged in mainline kernel
- ⇒ HUGE collection of patches
- ⇒ 80% already mainlined (timers, interrupt handlers, tracing infrastructure...)

Normal interrupt preempts the task and executes the handler function :



Threaded interrupt schedules the `thread_fn` function (*top/bottom half approach*) :



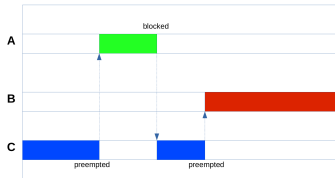
Forced threaded interrupts just acknowledges the device ($< 1 \mu s$) :



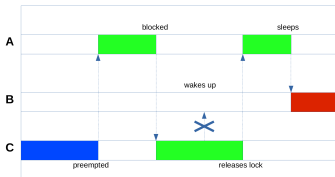
⇒ merged in mainline since 2009

⇒ all interrupts as threads `threadirqs` in mainline since 2011

Priority inversion :



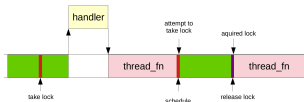
Priority inheritance :



Most of the spinlock are transformed in rt-mutexes. The others are now called raw_spinlocks.



There are in fact needed now with the threaded IRQs !



Usual way to tests real-time kernel behavior is to use the `rt-tests` suite :

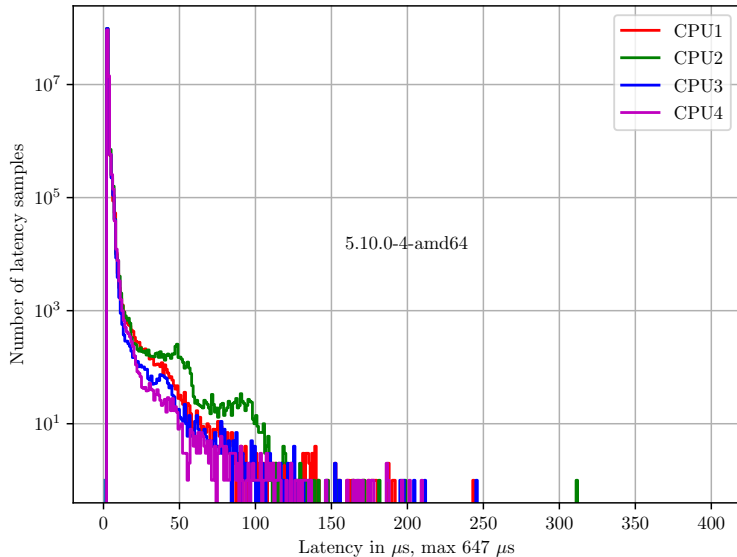
- timer latency
- signal latency
- functioning of priority-inheritance mutexes
- ...

The main program is `cyclictest` :

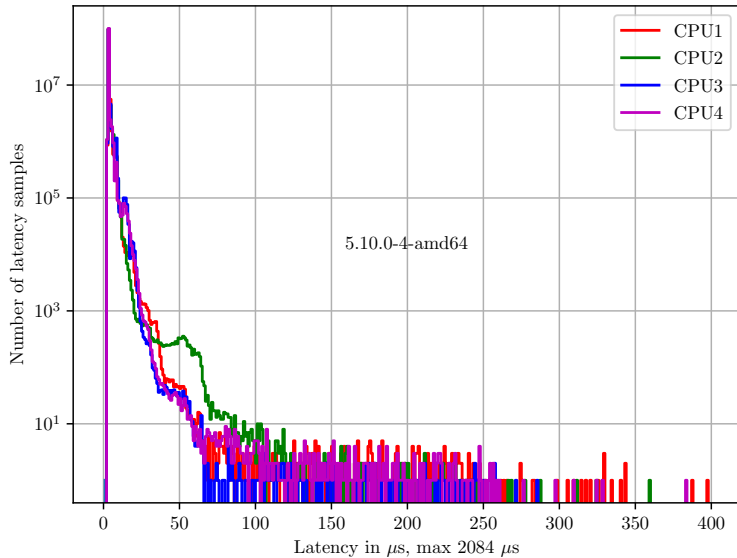
- `-p, --prio=PRIO` of first thread (default 80)
- `-i, --interval=INTV` of first thread (default $1000\mu\text{s}$)
- `-l, --loops=LOOPS` of first thread (default 0=endless)
- `-D, --duration=TIME`
- `-t, --threads=NUM` (default 1, empty means #CPU)
- `-m, --mlockall` lock current and future memory alloc
- `-a, --affinity=PROC-SET`
- `-S, --smp = -t -a`
- `-h, --histogram=US` max latency tracked

⇒ `cyclictest -D6h -m -S -p95 -i200 -h400`

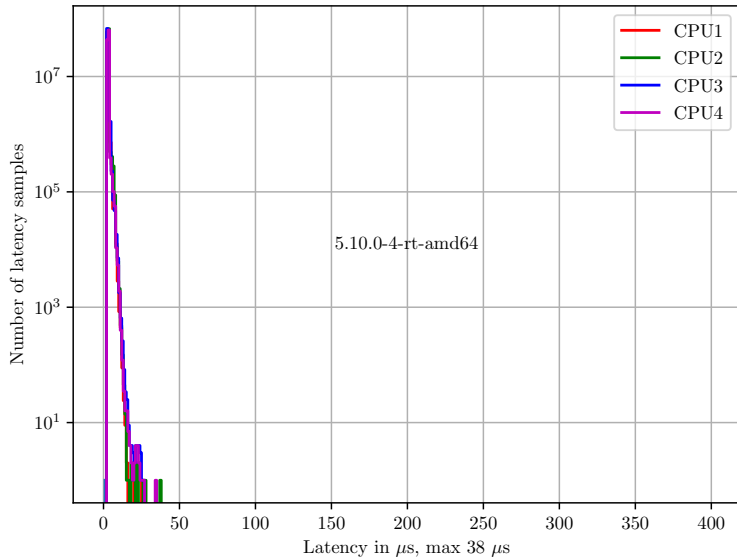
cyclictest -D6h -m -S -p 95 -i200 -h400



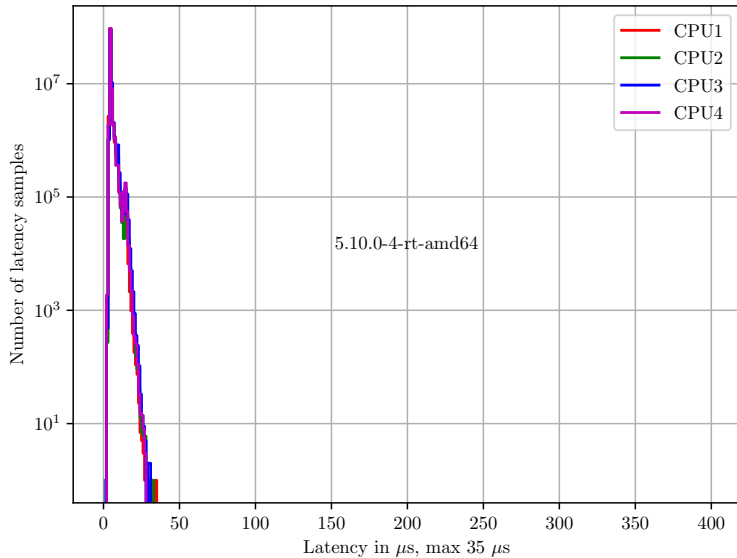
cyclictest -D6h -m -S -p 95 -i200 -h400



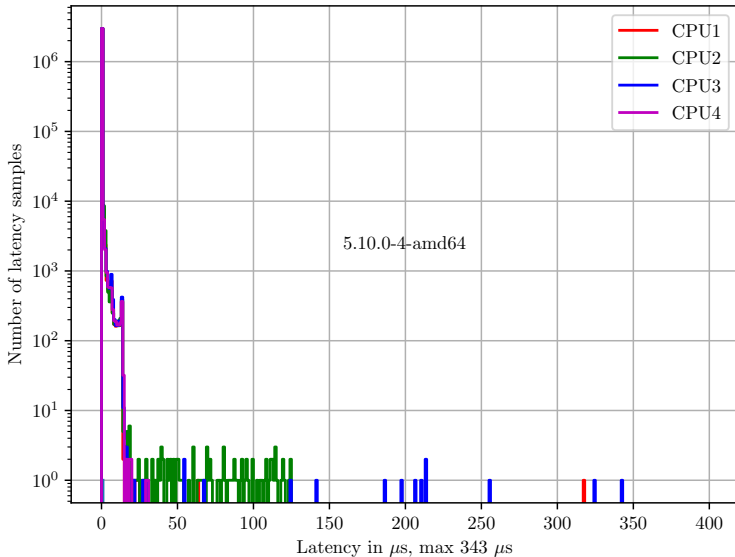
cyclictest -D6h -m -S -p 95 -i200 -h400



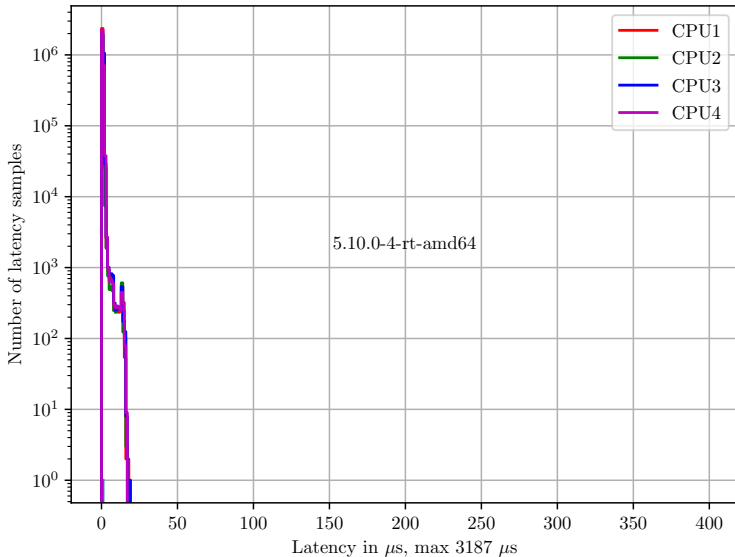
cyclictest -D6h -m -S -p 95 -i200 -h400



cyclictest -D10m -m -S -p 95 -i200 -h400



cyclictest -D10m -m -S -p 95 -i200 -h400



Oliveira, Casini, Oliveira & Cucinotta papers :

- model the IRQs, NMIs, thread synchronization mechanisms into an automata
 - model stable between Linux versions (contrary to a model entirely generated from traces)
 - able to find some linux code errors
 - apply in to get a model of the Linux+PREEMPT_RT patches latency
- ⇒ not just the latency but root causes of the latency
- ⇒ ~ Oliveira 2020 PhD thesis
- automata compiled as a kernel module
- ⇒ new tracing infrastructure to understand the root causes of this latency
- ⇒ alternative to `cyclictest` and kernel `ftrace`, user-space `trace-cmd` and graphic interface `kernelshark`

2) Model of IRQs, NMIs and thread synchronization mechanisms

Journal of Systems Architecture Volume 107, August 2020, 101729

A Thread Synchronization Model for the PREEMPT_RT Linux
Kernel

Daniel B. de Oliveira^{a,b,c}, Rômulo S. de Oliveira^b, Tommaso Cucinotta^c

^a*RHEL Platform/Real-time Team, Red Hat, Inc., Pisa, Italy.*

^b*Department of Systems Automation, UFSC, Florianópolis, Brazil.*

^c*RETIS Lab, Scuola Superiore Sant'Anna, Pisa, Italy.*

Synchronization mechanisms between :

- IRQs (Interrupt ReQuests)
- NMI (Non Maskable Interrupts) = cannot ignore interrupts = hw errors, parity/ECC errors, system
- thread scheduling, context switch, ...
- locking : mutex, rwlocks, semaphores

Model with an automata $G = (X, E, F, \Gamma, x_0, X_m)$:

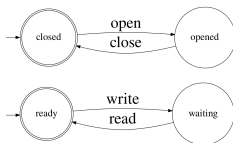
- X finite set of states
- E finite set of events
- $F : X \times E \rightarrow X$ transition function
- $\Gamma(x)$ set of events e such that $F(x, e)$ is defined in state x
- x_0 initial state
- X_m set of final states

On all the interesting Linux mechanisms ?

Modular approach relying on automata theory :

- generators, specifications, parallel composition
- computed automatically thanks to a dedicated software
Supremica IDE

Two generators or sub-automatas :

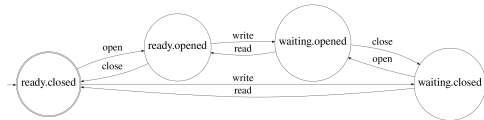


For a model of Linux synchronization mechanisms, this will be the minimal operations in kernel synchronization.

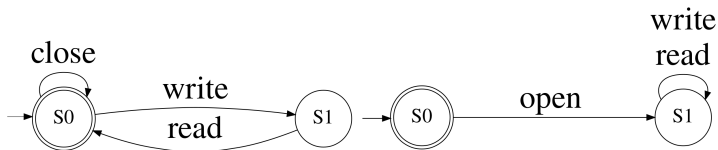
Parallel composition or synchronous composition of two generators $G_1 = (X_1, E_1, f_1, \Gamma_1, x_01, X_{m1})$ $G_2 = (X_2, E_2, f_2, \Gamma_2, x_02, X_{m2})$ introduces the notion of *private* events and *common* events

$$G_1 \parallel G_2 := Ac(X_1 \times X_2, E_1 \cup E_2, f_{1|2}, \Gamma_{1|2}, (x_01, x_02), X_{m1} \times X_{m2})$$

$$f_{1|2}((x_1, x_2), e) := \begin{cases} (f_1(x_1, e), f_2(x_2, e)) & \text{if } e \in \Gamma_1(x_1) \cap \Gamma_2(x_2) \\ (f_1(x_1, e), x_2) & \text{if } e \in \Gamma_1(x_1) \setminus E_2 \\ (x_1, f_2(x_2, e)) & \text{if } e \in \Gamma_2(x_2) \setminus E_1 \\ \text{undefined} & \text{otherwise} \end{cases}$$
$$\Gamma_{1|2}((x_1, x_2)) := \Gamma_1(x_1) \cup \Gamma_2(x_2)$$



Specifications between generators are also implemented as an automata.



For a model of Linux synchronization mechanisms, this will be for example the necessary conditions to call the scheduler.

And Supremica IDE computes the final automaton

- verifies there is no dead-lock
- final automata deterministic (only 1 transition to next step)

⇒ and Linux PREEMPT_RT too if the model is correct

Is the model really modeling the Linux kernel ?

Kernel event	Automaton event	Description
hw_local_irq_disable	preemptirq:irq_disable	Begin IRQ handler
hw_local_irq_enable	preemptirq:irq_enable	Return IRQ handler
local_irq_disable	preemptirq:irq_disable	Mask IRQs
local_irq_enable	preemptirq:local_irq_enable	Unmask IRQs
nmi_entry	irq_vectors:nmi	Begin NMI handler
nmi_exit	irq_vectors:nmi	Return NMI Handler

Kernel event	Automaton event	Description
preempt_disable	preemptirq:preempt_disable	Disable preemption
preempt_enable	preemptirq:preempt_enable	Enable preemption
preempt_disable_sched	preemptirq:preempt_disable	Disable preemption to call the scheduler
preempt_enable_sched	preemptirq:preempt_enable	Enables preemption returning from the scheduler
schedule_entry	sched:sched_entry	Begin of the scheduler
schedule_exit	sched:sched_exit	Return of the scheduler
sched_need_resched	sched:set_need_resched	Set need resched
sched_waking	sched:sched_waking	Activation of a thread
sched_set_state_runnable	sched:sched_set_state	Thread is runnable
sched_set_state_sleepable	sched:sched_set_state	Thread can go to sleepable
sched_switch_in	sched:sched_switch	Switch in of the thread under analysis
sched_switch_suspend	sched:sched_switch	Switch out due to a suspension of the thread under analysis
sched_switch_preempt	sched:sched_switch	Switch out due to a preemption of the thread under analysis
sched_switch_blocking	sched:sched_switch	Switch out due to a blocking of the thread under analysis
sched_switch_in_o	sched:sched_switch	Switch in of another thread
sched_switch_out_o	sched:sched_switch	Switch out of another thread

Kernel event	Automaton event	Description
mutex_lock	lock:rt_mutex_lock	Requested a RT Mutex
mutex_blocked	lock:rt_mutex_block	Blocked in a RT Mutex
mutex_acquired	lock:rt_mutex_acquired	Acquired a RT Mutex
mutex_abandon	lock:rt_mutex_abandon	Abandoned the request of a RT Mutex
write_lock	lock:rwlck_lock	Requested a R/W Lock or Sem as writer
write_blocked	lock:rwlck_block	Blocked in a R/W Lock or Sem as writer
write_acquired	lock:rwlck_acquired	Acquired a R/W Lock or Sem as writer
write_abandon	lock:rwlck_abandon	Abandoned a R/W Lock or Sem as writer
read_lock	lock:rwlck_lock	Requested a R/W Lock or Sem as reader
read_blocked	lock:rwlck_block	Blocked in a R/W Lock or Sem as reader
read_acquired	lock:rwlck_acquired	Acquired a R/W Lock or Sem as reader
read_abandon	lock:rwlck_abandon	Abandon a R/W Lock or Sem as reader

Name	States	Events	Transitions
<i>G01</i> Sleepable or runnable	2	3	3
<i>G02</i> Context switch	2	4	4
<i>G03</i> Context switch other thread	2	2	2
<i>G04</i> Scheduling context	2	2	2
<i>G05</i> Need resched	1	1	1
<i>G06</i> Preempt disable	3	4	4
<i>G07</i> IRQ Masking	2	2	2
<i>G08</i> IRQ handling	2	2	2
<i>G09</i> NMI	2	2	2
<i>G10</i> Mutex	3	4	6
<i>G11</i> Write lock	3	4	6
<i>G12</i> Read lock	3	4	6
<i>S01</i> Sched in after wakeup	2	5	6
<i>S02</i> Resched and wakeup sufficiency	3	10	18
<i>S03</i> Scheduler with preempt disable	2	4	4
<i>S04</i> Scheduler doesn't enable preemption	2	6	6
<i>S05</i> Scheduler with interrupt enabled	2	4	4
<i>S06</i> Switch out then in	2	20	20
<i>S07</i> Switch with preempt/irq disabled	3	10	14
<i>S08</i> Switch while scheduling	2	8	8
<i>S09</i> Schedule always switch	3	6	6
<i>S10</i> Preempt disable to sched	2	3	4
<i>S11</i> No wakeup right before switch	3	5	8
<i>S12</i> IRQ context disable events	2	27	27
<i>S13</i> NMI blocks all events	2	34	34
<i>S14</i> Set sleepable while running	2	6	6
<i>S15</i> Don't set runnable when scheduling	2	4	4
<i>S16</i> Scheduling context operations	2	3	3
<i>S17</i> IRQ disabled	3	4	4
<i>S18</i> Schedule necessary and sufficient	8	9	27
<i>S19</i> Need resched forces scheduling	7	25	53
<i>S20</i> Lock while running	2	16	16
<i>S21</i> Lock while preemptive	2	16	16
<i>S22</i> Lock while interruptible	2	16	16
<i>S23</i> No suspension in lock algorithms	3	10	19
<i>S24</i> Sched blocking if blocks	3	10	20
<i>S25</i> Need resched blocks lock ops	2	15	17
<i>S26</i> Lock either read or write	3	6	6
<i>S27</i> Mutex doesn't use rw lock	2	11	11
<i>S28</i> RW lock does not sched unless block	4	11	22
<i>S29</i> Mutex does not sched unless block	4	7	16
<i>S30</i> Disable IRQ in sched implies switch	5	6	10
<i>S31</i> Need resched preempts unless sched	3	5	12
<i>S32</i> Does not suspend in mutex	3	5	11
<i>S33</i> Does not suspend in rw lock	3	8	16
Model	9017	34	20103


The PREEMPT RT task model has:

- 12 generators, 33 specifications, 9017 states, 23103 transitions
- deterministic automata


3) Model the latency bound

Proceedings 32nd Euromicro Conference on Real-time Systems (ECRTS 2020)


Demystifying the Real-Time Linux Scheduling Latency

Daniel Bristot de Oliveira 


Red Hat, Italy
bristot@redhat.com

Daniel Casini 

Scuola Superiore Sant'Anna, Italy
daniel.casini@santannapisa.it

Rômulo Silva de Oliveira 

Universidade Federal de Santa Catarina, Brazil
romulo.deoliveira@ufsc.br

Tommaso Cucinotta 

Scuola Superiore Sant'Anna, Italy
tommaso.cucinotta@santannapisa.it

Idea : taking a subset of the model (9/12 generators, 14/33 specifications) and perform a case analysis to model the latency

Thread scheduling latency is actual start F (after context switch) - expected activation of highest priority A :

► **Definition 1** (Thread Scheduling Latency). *The scheduling latency experienced by an arbitrary thread $\tau_i^{THD} \in \Gamma^{THD}$ is the longest time elapsed between the time A in which any job of τ_i^{THD} becomes ready and with the highest priority, and the time F in which the scheduler returns and allows τ_i^{THD} to execute its code, in any possible schedule in which τ_i^{THD} is not preempted by any other thread in the interval $[A, F]$.*

Model divided in two parts :

- blocking/scheduling
- interference from IRQ/NMI

$$L = L^{IF} + I^{NMI}(L) + I^{IRQ}(L).$$

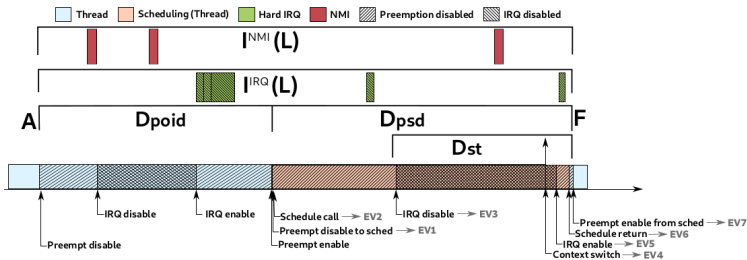
$I^{NMI}(L)$ and $I^{IRQ}(L)$ can be estimated from experiments and tests
 L^{IF} is the context switch and priority inversion blocking \Rightarrow try to find a bound

- EV1** The necessary conditions to call the scheduler need to be fulfilled: IRQs are enabled, and preemption is disabled to call the scheduler. It follows from rule R5 and R6;
- EV2** The scheduler is called. It follows from R12;
- EV3** In the scheduler code, IRQs are disabled to perform a context switch. It follows from rule R8;
- EV4** The context switch occurs. It follows from rule R13 and R14;
- EV5** Interrupts are enabled by the scheduler. It follows from R5;
- EV6** The scheduler returns;
- EV7** The preemption is enabled, returning the thread its own execution flow.

Mutually case exclusive analysis :

- i-a** if RHP_i occurs between events EV1 and EV2, i.e., after that preemption has been disabled to call the scheduler and before the actual scheduler call (black in Figure 21);
 - i-b** if RHP_i occurs in the scheduler between EV2 and EV3, i.e., after that the scheduler has already been called and before interrupts have been disabled to cause the context switch (pink in Figure 21);
 - i-c** if RHP_i occurs in the scheduler between EV3 and EV7, i.e., after interrupts have already been masked in the scheduler code and when the scheduler returns (brown in Figure 21);
- In case (ii), RHP_i occurred when the current thread $\tau_j^{THD} \in \Gamma_{LP_i}^{THD}$ is not in the scheduler execution flow. Based on the automaton of Figure 21, two sub-cases are further differentiated:
- ii-a** when RHP_i is caused by an IRQ, and the currently executing thread may delay RHP_i only by disabling interruptions (green in Figure 21).
 - ii-b** otherwise (blue in Figure 21).

Param.	Length of the longest interval
D_{PSD}	in which preemptions are disabled to schedule.
D_{PAIE}	in which the system is in state <code>pe_ie</code> of Figure 21.
D_{POID}	in which the preemption is disabled to postpone the scheduler or IRQs are disabled.
D_{ST}	between two consecutive occurrences of EV3 and EV7.



■ Figure 22 Reference timeline.

$$L^{IF} \leq \max(D_{ST}, D_{POID}) + D_{PAIE} + D_{PSD},$$

$$L = \max(D_{ST}, D_{POID}) + D_{PAIE} + D_{PSD} + I^{NMI}(L) + I^{IRQ}(L)$$

4) Conclusion

PREEMPT_RT patches :

- several advantages on other Linux real-time approaches
- soon in main git repository (2021?)
- some optimizations already merged

Model with the automata approach :

- 12 generators, 33 specifications
- ⇒ model tractable
- ⇒ model stable between versions
- ⇒ can help to find bugs
- claims that Linux+patches is deterministic and the scheduling latency is bounded

Details not covered in this presentation that are worth to look at :

- extension of the perf tool and automata compiled inside a kernel module
- comparison with cyclicttest and ftrace