**Red Hat Enterprise Linux**

**Sant'Anna**
School of Advanced Studies – Pisa

**UNIVERSIDADE FEDERAL DE SANTA CATARINA**

# Demystifying the Real-Time Linux Scheduling Latency

**Daniel Bristot de Oliveira**, Daniel Casini, Rômulo Silva de Oliveira and Tommaso Cucinotta
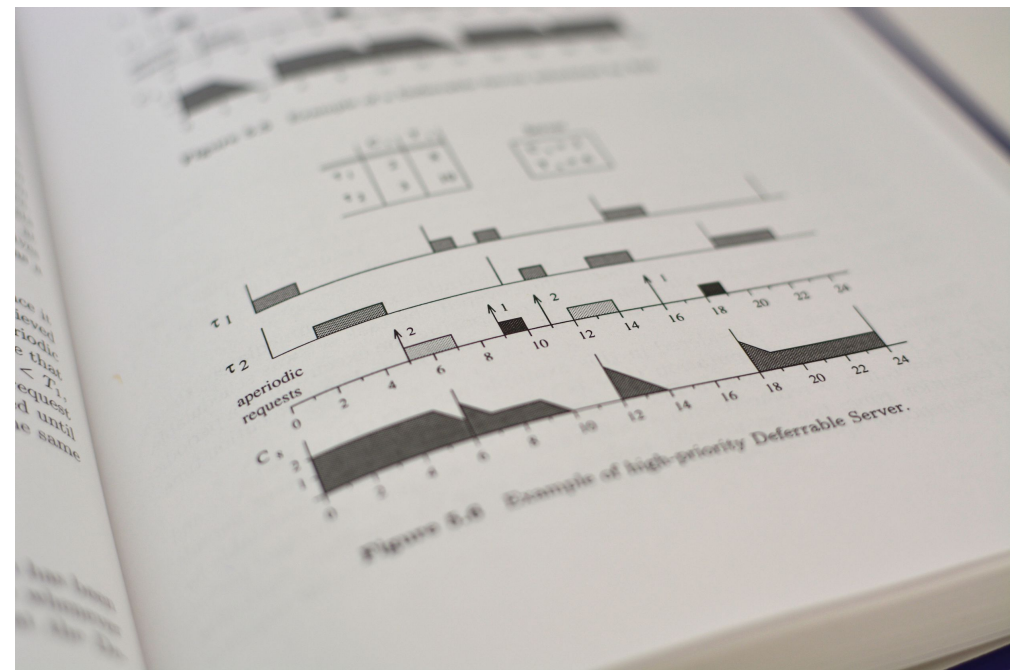**Principal Software Engineer**

**Red Hat**

# Real-Time Linux

Demystifying The Real-Time Linux Scheduling Latency – 32$^{nd}$ Euromicro Conference on Real-Time Systems – ECRTS'20

Red Hat

# "Real-Time" Linux

Demystifying The Real-Time Linux Scheduling Latency – 32$^{nd}$ Euromicro Conference on Real-Time Systems - ECRTS'20

# Real–Time Linux vs Real–Time theory

## Experimental vs Analytical

Demystifying The Real-Time Linux Scheduling Latency – 32nd Euromicro Conference on Real-Time Systems - ECRTS'20

# Real-Time Linux vs Real-Time theory

## Real-time analysis



- Based on the timing description of the system
- Capture all behaviors
- Precisely define the worst cases
- But depends on a precise definition of the system

Demystifying The Real-Time Linux Scheduling Latency – 32nd Euromicro Conference on Real-Time Systems - ECRTS'20

# Real-Time Linux vs Real-Time theory

## Linux approach



- Linux was adapted to become a RTOS

- PREEMPT_RT: *De facto* standard

- Evaluated (mainly) with cyclictest

- Cyclictest:
  - Practical: lightweight and out-of-the-box
  - It is a "black-box" test
  - No demonstration
  - Does not provide evidence of "root-cause"

# Why don't we apply RT analysis on Linux?

# Linux is complex



- Lots of contexts

- Lots of hacks

- Lots of information

- Fast pacing

- ...

# The PREEMPT_RT thread model



It defines the *specifications* of threads synchronization:



Figure 22: *S08* Switch while scheduling.

Figure 23: *S03* Scheduler called with preemption disabled.

Figure 20: *S05* Scheduler called with interrupts enabled.

Figure 21: *S07* Switch with interrupts and preempt disabled.

# Demystifying the Real-Time Linux Scheduling Latency

## Approach

### Formal specification



### Scheduling latency bound



### Measurement and analysis

# From formal specification to synchronization rules

## Formally backed natural language arguments



- Generators
  - Basic/Independent behavior
  - e.g., irq_disable/enable, scheduler call
- Translated into a set of operations

- Specifications
  - Relations among generators
  - e.g., necessary conditions to call the scheduler
- Translated into a set of synchronization rules

# Scheduling latency definition

The **scheduling latency** experienced by an arbitrary thread τ is:

- the **longest time** elapsed **between** the *time A* in which any job of τ becomes **ready and with the highest priority**,
- and the *time F* in which the scheduler returns and allows τ to **execute its code**.

From the first necessary condition to *set need resched*, to the the last action after the scheduling, which is *enabling preemption* after the return from __schedule().

# Interference and blocking

```
void __sched notrace __schedule(bool

struct task_struct *prev, *next;
unsigned long *switch_count;
struct rq_flags rf;
struct rq *rq;
int cpu;

cpu = smp_processor_id();
rq = cpu_rq(cpu);
prev = rq->curr;

schedule_debug(prev, preempt);

if (sched_feat(HRTICK))
        hrtick_clear(rq);

local_irq_disable();
rcu_note_context_switch(preempt);
```

The **scheduling latency** is caused by:

● **Blocking** from the current (and so lower) priority thread;
  ● Including scheduling.

● **Interference** from IRQs and NMI.

The scheduling latency in this paper refers to the delay between the notification of a new highest priority thread, to point in which this thread starts running its own code.

The highest priority thread can belong to any scheduler: the analysis is scheduler independent.

**Red Hat**

# Blocking bound

## From the specification that bounds the block to a timeline

# Timeline and cases

## All possible cases



i-a/i-b

i-c

ii-a/ii-b

IRQ disable    IRQ enable    Schedule call → **EV2**    IRQ disable → **EV3**

Preempt disable to sched → **EV1**

Preempt enable

Preempt disable

Preempt enable from sched → **EV7**

Schedule return → **EV6**

IRQ enable → **EV5**

Context switch → **EV4**

# Blocking variables

- **D<sub>POID</sub>**: preemption or interrupts disabled to postpone the scheduler;

- **D<sub>PAIE</sub>**: preemption and interrupts enabled, as a transient state from **poid** to **psd**; when scheduling a new highest priority thread.

- **D<sub>PSD</sub>**: preemption disable to schedule;

- **D<sub>ST</sub>**: delay caused by the scheduling tail; the "non return" point in which a new arrived task will have to wait for the current scheduling operation to finish before scheduling.

In the model, the preemption control is specialized into two different operations: to *postpone the scheduler* (the most known behavior) or to *protect the execution of the __schedule()* function from recursion.

Demystifying The Real-Time Linux Scheduling Latency - 32<sup>nd</sup> Euromicro Conference on Real-Time Systems - ECRTS'20

# Timeline and cases

## Variables in the the timeline

Demystifying The Real-Time Linux Scheduling Latency – 32$^{nd}$ Euromicro Conference on Real-Time Systems - ECRTS'20

# Timeline and cases

## IRQ and NMI interference

# And the *scheduling latency* bounds to:



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

The bound considers all possible cases. Note that the Latency *L* is present in both sides of the equation.

So, L is bounded by the least positive value fulfilling the equation (like on RTA).

# Interrupts are workload dependent

- Instead of proposing "the best" interrupt characterization, the rtsl reports the scheduling latency based on some well-known characterizations:
  - No interrupt
  - Worst single interrupt
  - Single occurence of all interrupts
  - Sporadic
  - Sliding window (Author's preferred)
  - Sliding window with oWCET

This topic was heavily discussed at the Real-time Micro Conference (inside Linux Plumbers) in 2019, more info here:

Red Hat

# A practical scheduling latency estimation tool

## Method and challenges



- Based on the latency bound
- The latency bound is based on the model
- The *model* is based on tracing of events
  - but high frequency events
    - hundreds MB/sec/CPU
- Challenges:
  - To minimize the (runtime) overhead
  - Work out-of-the-box

# rt_sched_latency (rtsl)



Based on **perf**

Works in two phases:

- The **record** mode saves the trace data;

- The **report** mode process the trace and does the analysis.

# record phase

## Low overhead trace recording



Kernel

tracepoints

latency parser

perf buffer

perf script record rtsl

perf.data

- Filters the high frequency trace
  - Doing in-kernel processing
- For blocking variables:
  - Reports only the discover of new max values
- For IRQ and NMI:
  - Reports one event for each occurrence
- Discounts the interference:
  - e.g., IRQ interference on a **poid**

# report phase

## Low overhead trace recording

```
┌──────────────┐
│  perf script │ ──→  Analysis
│  report rtsl │
│              │
│   perf.data  │──→  Chart
└──────────────┘
```

- After the capture, analyzes the trace.
  - All in user-space.
- Most of the analysis is done in python
  - Easy to extend
- Two outputs:
  - Textual: good for debug
  - Chart: good comparisons (and papers :-))
- Does a per-cpu scheduling latency analysis
  - Using different IRQ/NMI characterization...

# rtsl report output

## Textual output

```
Interference Free Latency:
    paie is lower than 1 us -> neglectable
    latency = max(poid,    dst) + paie +    psd
       42212 = max(22510, 19312) +    0 + 19702
Cyclictest:
    Latency =      27000 with Cyclictest
No Interrupts:
    Latency =      42212 with No Interrupts
Sporadic:
    INT:       oWCET              oMIAT
    NMI:           0                  0
     33:       16914             257130
     35:       12913               1843 <- oWCET > oMIAT
    236:       20728               1558 <- oWCET > oMIAT
    246:        3299            1910321
    Did not converge.
```

```
continuing....
Sliding window:
    Window: 42212
            NMI:           0
             33:       16914
             35:       14588
            236:       20728
            246:        3299
    Window: 97741
            236:       21029 <- new!
    Window: 98042
    Converged!
    Latency =      98042 with Sliding Window
```

# rtsl report output

## Chart output

Demystifying The Real-Time Linux Scheduling Latency – 32$^{nd}$ Euromicro Conference on Real-Time Systems – ECRTS'20

# Experiments

The experiments passed
by the artifact evaluation!

- Scheduling latency measurements on two systems:
  - workstation: eighth CPUs
  - server: twelve CPUs server
- Experiments:
  - Single-core
    - Different duration
    - Different workload
  - Multi-core
- Running in parallel with cyclictest
- Note: The goal of the experiments is to demonstrate the tool, not to define worst values.

Demystifying The Real-Time Linux Scheduling Latency – 32nd Euromicro Conference on Real-Time Systems - ECRTS'20

# Single-core experiments



Legend: Cyclictest | No Interrupts | Worst Single Interrupt | Single (Worst) of Each Interrupt | Sliding Window | Sliding Window with oWCET

1.a) Idle

1.b) CPU Intensive

1.c) I/O Intensive

2.a) 15 min.

2.b) 60 min.

2.c) 180 min.

Demystifying The Real-Time Linux Scheduling Latency – 32nd Euromicro Conference on Real-Time Systems – ECRTS'20

# Multicore experiments



Legend: Cyclictest · No Interrupts · Worst Single Interrupt · Single (Worst) of Each Interrupt · Sliding Window · Sliding Window with oWCET

3.a) Workstation Idle

3.b) Workstation CPU Intensive

3.c) Workstation I/O Intensive

4.a) Server I/O Intensive

# Conclusions

For more information about this paper, like source code, other comments, Q&A, check its companion page!

- The PREEMPT_RT preemption model is deterministic, and the scheduling latency is bounded.
- The approach presented in this paper opens the door for a new set of real-time analysis for Linux;
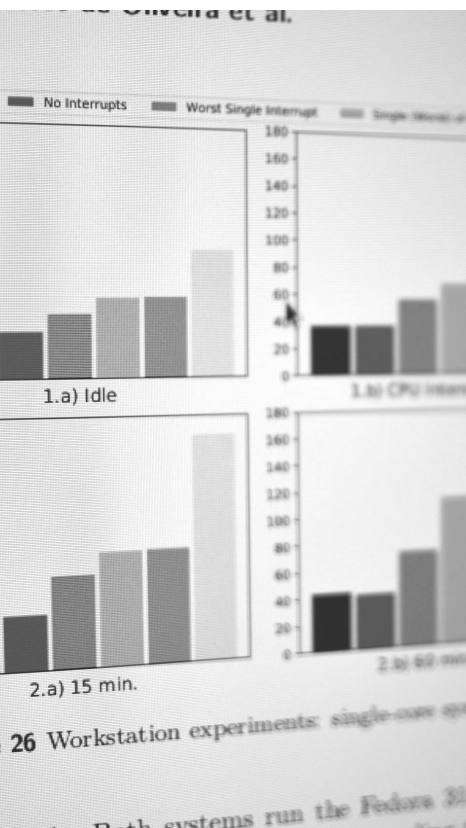  - The analytical interpretation of Linux thread model developed in this paper untight the Linux complexity, enabling the reasoning at a more sophisticated level.
- Even though rtsl finds higher scheduling latency values, they are still low enough to justify Linux as RTOS on the current scenarios.
- rtsl is practical, and resolves many problems of cyclictest.
  - E.g., it can be used to point to the root causes of the latency;
  - But still can, and should, be improved:
    - Both with code, and other analysis.

# Thank you

Red Hat is the world's leading provider of enterprise

open source software solutions. Award-winning

support, training, and consulting services make

Red Hat a trusted adviser to the Fortune 500.

in    linkedin.com/company/red-hat

▶    youtube.com/user/RedHatVideos

f    facebook.com/redhatinc

🐦    twitter.com/RedHat

**Red Hat**