





A theorem for the RT scheduling latency (and a measurement tool too)

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









Episode III - Showing the math

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



A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020



Real-Time Linux

A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020



4

"Real-Time" Linux



Real-Time Linux vs Real-Time theory

Experimental vs Analytical

(root@realtime-d	000) cyc : cat 01 ~l≢ cyclicto	estsmp -p 05	- 173		
# /dev/cpu_dma_1					
policy: fife: la	oadavg: 14.90	6.21 3.98 2/387	2735923		
T: 0 (2735898)	P:95 I:1000 C:	66520 Min:	4 Act:	5 Avg:	
T: 1 (2735899)	P:95 I:1500 C:	44341 Min:	4 Act:		
T: 2 (2735900)	P:95 I:2000 C:	33256 Min:	4 Act:	6 Avg:	
T: 3 (2735901)	P:95 I:2500 C:	26598 Min:	4 Act:		
To 4 (2735902)	P:95 I:3000 C:	22162 Min:	4 Act:	5 Avg:	
To 5 (2735903)	P:95 I:3500 C:	18903 Min:	4 Act:	6 Avg:	
T: 6 (2735904)	P:95 I:4000 C:	16617 Min:	4 Act:		
To 7 (2735905)	P:95 I:4500 C:	14769 Min:	4 Act:		
T & (2735906)	P:95 I:5000 C:	13290 Min:	4 Act:	6 Avg:	
T. 6 (3735907)	P:95 I:5500 C:	12080 Min:	4 ACL:	12 Avg:	
TALA (2735908)	P:95 I:6000 C:	11002 Min:	A Act:		
T-11 (2735909)	P:95 I:6500 C:	10219 Min:	5 Act:		
(2735910)	P:95 I:7000 C:	9488 Min:	5 Act:		
(2735911)	P:95 I:7500 C:	8850 Min:			
Tila (2735912)	P:95 I:8000 C:	8300 Min.			
(2735913)	P:95 I:8500 C:	7801 Min:			
7.16 (2735914)	P:95 I:9000 C:	6087 Min:			
T:17 (2735915)	P:95 I:9500 C:	6638 Min:	5 Act:		







Real-Time Linux vs Real-Time theory

Linux approach

revenue i semicul che i car .	uculctest.txt			
[root@realtime-01 ~]# cyclicte				
policy: fife: loadavg: 14.90 (5.21 3.98 2/387	2735923		
T: 0 (2735898) P:95 I:1000 C:	66520 Min:	4 Act:	5 Avg: 5 M	
T: 1 (2785899) P:95 I:1500 C:	44341 Min:	4 Act:		
T: 0 (2735000) P:95 I:2000 C:	33256 Min:	4 Act:	6 Avg: 5 M	
T: 3 (2735901) P:95 I:2500 C:	26598 Min:	4 Act:		
P. A (9735003) P:95 I:3000 C:	22162 Min:	4 Act:	5 Avg: 5 M	
(1735003) P:05 T:3500 C:	18993 Min:	4 Act:		
11 5 (2735503) P.05 T.4000 C:	16607 Min:	4 Act:		
T: 6 (2730904) P:93 1:4000 C:	14769 Min:	4 Act:		
T: 7 (2735905) P:95 1:4500 C:	13290 Min:	4 Act:		
T: 8 (2735906) P:95 1:5000 C.	12080 Min:	4 Act:		
T: 9 (2735907) P:95 1:5500 C.	11057 Min:	8 Act:	12 Avg: 13 M	
T:10 (2735908) P:95 I:6000 C.	10210 Min:	4 Act:		
T:11 (2735909) P:95 I:6500 C:	0488 Min:	5 Act:		
T:12 (2735910) P:95 I:7000 C:	0054 Min:	5 Act:		
T:13 (2735911) P:95 I:7500 C:	ezaA Min:			
T:14 (2735912) P:95 I:8000 C:	7007 Min:			
T:15 (2735913) P:95 I:8500 C:	7376 Min:			
TILE (2735914) P:95 I:9000 C:	6983 Min:			
TIT (2735915) P:95 I:9500 C:	6638 Min:	5 Act:		
	. 0000	e Act:		

- 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 "closed-box" test
 - No demonstration
 - Does not provide evidence of "root-cause"



7

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
- Often overly-simplified



g: 3 cat occlictest.txt # cyclictest --smp -p 95 -m hcy set to Ous /g: 14.90 6.21 3.98 2/387 2735923

I:1000	C:	66520	Min:		Act:	
I:1500	C:	44341	Min:		Act:	
I:2000	C:	33256	Min:	4	Act:	
I:2500	C:	26598	Min:	4	Act:	
I:3000	C:	22162	Min:	4	Act:	
I:3500	c:	18993	Min:	4	Act:	
I:4000	C:	16607	Min:	4	Act:	
I:4500	C:	14762	Min:	4	Act:	
I:5000	C:	13290	Min:	4	Act:	
T:5500	c:	12080	Min:	4	Act:	
T:6000	c:	11002	Min:		ACT:	
T:6500	с:	10219	Min:	4	ACT:	
T:7000	с:	9488	Min:		ACT:	
T:7500	с:	8854	Min:		Act:	
T:8000	с:	8200	Min:		Act:	
T:8500	с:	780I	Min:		Act:	
I:9000	С:	7370	Man:		Act:	
I:9500	С:	6983	Mill.			
- 1000	A C:	663	0 11111		: Act:	

But, I like both.





A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020

Introduction

cyclictest --smp -p 95 -m hcy set to Ous

I:1000	C:	66520	Min:		Act:	
I:1500	C:	44341	Min:		Act:	
I:2000	C:	33256	Min:	4	Act:	
I:2500	c:	26598	Min:	4	Act:	
I:3000	c:	22162	Min:	4	Act:	
I:3500	c:	18993	Min:	4	Act:	
I:4000	c:	16607	Min:	4	Act:	
T:4500	c:	14769	Min:	4	Act:	
T:5000	c:	13290	Min:	4	Act:	
T-5500	c:	12080	Min:		Act:	
T.6000	C:	11002	Min:		Act:	
T-6500	c:	10219	Min:	4	Act:	
T . 7000	c:	9488	Min:		ACT:	
T.7500	c:	8850	Min:		ACL	
T-8000	с:	8200	Min:		Act:	
T-8500	с:	780I	Min:		Act:	
T-9000	c:	7370	Min:		Act:	
T . 9500	C:	6983	Min:	- 5		
		663	g Min:			

Demystifying the Real-Time Linux Scheduling Latency

Daniel Bristot de Oliveira ⁽²⁾ 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

- Abstract -

Linux has become a viable operating system for many real-time workloads. However, the black-box approach adopted by cyclictest, the tool used to evaluate the main real-time metric of the kernel, the scheduling latency, along with the absence of a theoretically-sound description of the in-kernel behavior, sheds some doubts about Linux meriting the real-time adjective. Aiming at clarifying the PREEMPT_RT Linux scheduling latency, this paper leverages the *Thread Synchronization Model* of Linux to derive a set of properties and rules defining the Linux kernel behavior from a scheduling perspective. These rules are then leveraged to derive a sound bound to the scheduling latency, considering all the sources of delays occurring in all possible sequences of synchronization events in the kernel. This paper also presents a tracing method, efficient in time and memory overheads, to observe the kernel events needed to define the variables used in the analysis. This results in an easy-to-use tool for deriving reliable scheduling latency bounds that can be used in practice. Finally, an experimental analysis compares the cyclictest and the proposed tool, showing that the proposed method can find sound bounds faster with acceptable overheads.

2012 ACM Subject Classification Computer systems organization \rightarrow Real-time operating systems

Keywords and phrases Real-time operating systems, Linux kernel, PREEMPT_RT, Scheduling latency

Digital Object Identifier 10.4230/LIPIcs.ECRTS.2020.9

Supplementary Material ECRTS 2020 Artifact Evaluation approved artifact available at https://doi.org/10.4230/DARTS.6.1.3.

Supplement material and the code of the proposed tool is available at: https://bristot.me/ demystifying-the-real-time-linux-latency/

Funding This work has been partially supported by CAPES, The Brazilian Agency for Higher Education, project PrInt CAPES-UFSC "Automation 4.0."

Acknowledgements The authors would like to thank Thomas Gleixner, Peter Zijlstra, Steven Rostedt, Arnaldo Carvalho De Melo and Clark Williams for the fruitful discussions about the model, analysis, and tool.

© Daniel Bristot de Oliveira, Daniel Casini, Rômulo Silva de Oliveira, and Tommaso Cucinotta; licensed under Creative Commons License CC-BY 32nd Euromicro Conference on Real-Time Systems (ECRTS 2020). Editor: Marcus Völp; Article No.9; pp. 9:1–9:23 Leibniz International Proceedings in Informatics





📥 Red Hat

Lipits Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Episode I: getting formal



Math side: Talk is cheap...





Demystifying the Real-Time Linux Scheduling Latency

Approach





13

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 **T** is:

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

ready AND with the highest priority:

It covers the case in which these **two actions are not a single event.**

It is scheduler independent.

There is only one highest priority thread on a CPU: it is the one selected to run by the scheduler.

😓 Red Hat

Interference and blocking

void __sched notrace __schedule(b

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);

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.

This are well established terms in the real-time scheduling literature:

Interference from higher priority, blocking from lower priority.



Blocking bound

From the specification that bounds the block to a timeline







Blocking bound

Scheduling latency: start

- The **longest time** elapsed **between** the *time* A in which any job of **T** becomes **ready and with the highest priority**:
 - Generalized to the **need_resched** event
 - Works for all schedulers
 - cyclictest does not work for DL with NR_TASKS > CPUs.
 - Works for all conditions
 - E.g., a throttled DL task after a replenishment will cause a need resched without a wakeup.
 - Has preempt and IRQ disabled as necessary conditions
 - So we use the occurrence of the first **necessary condition** as the starting point of the critical window.
 - E.g., when preemption was disabled for the first time.

The wakeup is the only event that causes a need resched, and that is why it was not used here.

But ready means that the task was awakened.



Blocking bound

Scheduling latency: end

- And the *time F* in which the scheduler returns and allows T to **execute** its code.
 - Generalized to the **preempt_enable** after **__schedule()**
 - Implies that the system **crossed the context switch** code path.
 - Context switch implies __schedule()
 - Context switch needs:
 - Preempt disable to schedule as necessary condition
 - irqs disabled by thread as necessary condition

We are looking for a safe-bound, and so we have to put pessimism values.

We can latter reduce the pessimism, but with safe arguments.



How do we bound that?







Timeline and cases

All possible cases



Timeline and cases

Variables in the the timeline







Blocking variables

- **Dpoid**: preemption or interrupts disabled to postpone the scheduler;
- DPAIE: preemption and interrupts enabled, as a transient state from **poid** to **psd**; when scheduling a new highest priority thread.
- **Dpsp**: preemption disable to schedule;
- DsT: 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.



Timeline and cases

Variables in the the timeline





Timeline and cases

IRQ and NMI interference



And the scheduling latency bounds to:

By leveraging the individual bounds on L^{17} is overall bound that is valid for all the Lemma 7.

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

The lemma follows by noting that case (he clusive and cover all the possible sequences of meed_reached, to the time instant in which Definition 1), and the right-hand side of Equations 2, 3, 4, and 5.

Theorem 8 summarizes the results derived is u Theorem 8. The scheduling latency experiment is a positive value that fulfills the following of $L = max(D_{ST}, D_{POID}) + D_{PAIE} + D_{PSD} + I^{**}$ and The theorem follows directly from Lemma is a scale to the IRQ handler, the present

$$L = max(Dst, Dpoid) + Dpaie + Dpsd + I^{NMI}(L) + I^{IRQ}(L)$$

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).



Timeline and cases

IRQ and NMI interference



L = max(Dst, Dpoid) + Dpaie + Dpsd + $I^{NMI}(L) + I^{IRQ}(L)$



And the scheduling latency bounds to:

By leveraging the individual bounds on L^{17} is overall bound that is valid for all the Lemma 7.

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

The lemma follows by noting that cases (in clusive and cover all the possible sequences of meed_reached, to the time instant in which Definition 1), and the right-hand side of Equations 2, 3, 4, and 5.

Theorem 8 summarizes the results derived is u Theorem 8. The scheduling latency experiment (nest positive value that fulfills the following in $L = max(D_{ST}, D_{POID}) + D_{PAIE} + D_{PSD} + I^{**}$ and The theorem follows directly from Lemma in the present to the IRQ handler, the present

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:





Episode II: getting practical (and efficient)



A practical scheduling latency estimation tool

Method and challenges

 	na ciunt e	Ac + ca	10 mg	CLICLE	st.txt					
			ictes		1p -p 95					
		icy set	to 0	us						
		14.0		21 2 0	18 2/287	2725	:072			
						2100				
		I:1000	c:	66520	Min:		Act:		Avg:	
		I:1500	C:	44341	Min:		Act:		Avg:	
	P:95	I:2000	C:	33256	Min:		Act:		Avg:	
	Dens.	T:2500	C:	26598	Min:		Act:		Avg:	
		T . 2000	c .	22162	Min:		Act:		Avg:	
		1.3000	č.	18993	Min:		Act:		Avg:	
) 11:23	1:3000	·	16600	Min	4	Act:		Avg:	
6 (2735904		I:4000	C:	10000	Mint	4	Act:			
) P:95	I:4500	c:	14/02	Hin.	4	Act:		Avg:	
		I:5000	C:	13288	PIIII		Act:			
		I:5500	C:	12080	Min:		Act:	12		
) P:95	I:6000	C:	11082	Min:		Act:			
) P:95	I:6500	С:	10219	Min:		Act:			
	N P.95	I:7000	C:	9488	Min:					
[5 (5129ate	D-05	T:7500	C:	8850	Min:					
13 (513281)	1) 5:39	T-8000	c:	8200	Min:					
14 (273591)	7) L:23	T. 2500	c:	780I	Min:					
	3) P:90	1.0000	C:	7370	Min:					
	4) P:95	1:9000		6983	Min:					
	5) P:95	1:3200		663	g Min:					

• 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



Rtsl: a measurement tool



Kernel space:

- Rtsl events

User space:

- Rtsl command Python

Has three commands:

- The **record** command saves the trace data;
- The **report** command process the trace and does the analysis.
- The stats command produces a histogram of the thread variables

🧕 Red Hat

rtsl events

Low overhead tracer



- Hooks to events
 - Filters the high frequency trace
 - Doing in-kernel processing
 - Use a knob on debugfs to enable the tracing
- For blocking variables:
 - Reports all values or 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**



A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020

Experiments



Kernel changes

• The rtsl events depends on:

preemptirq tracepoints

- So it needs a "debug/trace" kernel (yeah...)
- But life finds a way
- Annotations on the preempt_disable to sched
 - No functional changes
- NMI tracepoints

• Or change in the current one to the extreme points of the handler

The parser was developed as a kernel module. In this way I can leave it off tree... but it would be better to have it in.

If we get it in, we can change the debugfs for the tracefs.



rtsl record

Trace recording





A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020

rtsl report

Trace processing



- Analyzes the trace!
 - All in user-space
- Most of the tool is done in python
 - Easy to extend the analysis (researchers like)
- Parses the trace file in parallel
 - Per cpu trace parsing (e.g., perf script-c \$...)
 - Generates per-cpu database with the data
 - In the rtsl_data/ dir
 - Uses a C trace-plugin create the database
- Database in a sqlite3 file



37

rtsl report

Data processing



- The analysis is done on the database
 - IRQ analysis needs to read data back and forth
 - Trace can reach tens of GB/per-cpu
- The analysis is done in parallel
- Two outputs:
 - Textual output
 - Charts
 - Using matplotlib



rtsl report output

Textual output

Interference 1	Free Latency	y:						
<pre>paie is lo latency = 42212 = Cyclictest: Latency = No Interrupts Latency =</pre>	ower than 1 max(poid, max(22510, 27000 t :	us -> negl dst) + p 19312) + with Cyclic	ectable aie + ps 0 + 1970 test	d 2	continuing Sliding window: Window: 42212 NMI: 33: 35:	0 16914 14588		
Sporadic:	42212 \	WICH NO INC	errupts		236:	20728		
INT: 0 NMI: 33: 35: 236: 2	oWCET 0 16914 12913 20728	oMIAT 0 257130 1843 1558	<- oWCET > <- oWCET >	oMIAT oMIAT	246: Window: 97741 236: Window: 98042 Converged! Latency =	21029 < 98042 w:	<- new! ith Sliding	Window
246: Did not co	3299 onverge.	1910321					0	

rtsl report output

Chart output



Red Hat



40

rtsl stats





• Monitor the thread variables

- poid/psd/dst/paie...
- Uses BCC
 - Saves histograms in kernel
 - Display in user-space
 - Can plot data



Experiments



Experiments

- 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.

The experiments passed by the artifact evaluation!





Single-core experiments





A theorem for the real-time scheduling latency (and a measurement tool too) - LPC 2020

43

Multicore experiments





Universe et a

Remarks

- The PREEMPT_R
 the scheduling late
 The approach press
 new set of real-time
 The analytical is
 in this paper units
- 2.a) 15 min. 26 Workstation experiments: single-converse 1. Both systems run the Federe 3.

The PREEMPT_RT preemption model is deterministic, and the scheduling latency is bounded.

- The approach presented in the 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.

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



44

1.a) Idle



45

rtsl vs cyclictest? nah

- They help the same people
 - But they do different things
- rtls is a more specific tool
 - Covers a single aspect: sched latency
 - Covers all cases at synchronization level
 - In the worst condition, even those that happened at different points in time.
 - With strong arguments
 - Depends on kernel features (PREEMPT_RT/preemptirq...)
- cyclictest is a more generic tool
 - Covers many aspects: external activation of the timer
 - Hardware delays? Hardware bugs?
 - Without analysis a closed-box test
 - Run on the potato that runs Linux
- rtsl adds only 4-ish us of overhead on cyclictest

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





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
 - facebook.com/redhatinc

twitter.com/RedHat



In the next episode....

