

The Linux Real-Time Task Model Extractor

31st IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS) May 9, 2025

<u>Björn Brandenburg</u>

MPI-SWS, Germany

Cédric Courtaud

Huawei Technologies, France

University of Southampton, UK

[Authors listed in alphabetical order. Work carried out while affiliated with MPI-SWS.]



MAX PLANCK INSTITUTE FOR SOFTWARE SYSTEMS

Filip Marković

Bite Ye *MPI-SWS, Germany*





European Research Council Established by the European Commission



LiME in a nutshell:

.....no source code & no static analysis of binary

blackbox threads running on Linux

online (= constant-memory) inference & conservative w.r.t. observations

well-studied in the RT literature & schedulability analysis available



······ no user input & no RT or kernel expertise required

The Problem Being Solved

Enter event or section title



The System-Model-Analysis-Prediction Cycle



Linux has Become a Serious RT Platform



The "Dual Expertise" Barrier



 $\tau_l = (C_l, T_l) \quad \tau_i = ($ $\tau_j = (C_j, \alpha_j(\Delta))$...

real-time task

State of the Art

→ *No tool support whatsoever!*

Doing it the hard way:

- **1**. Become familiar with Linux's low-level tracing facilities (*eBPF*).
- **2**. Become familiar with Linux system call semantics below libc level.
 - → *Dozens of relevant syscalls and flags...*
- **3**. Become familiar with applicable RT theory.
 - → *In particular, identify suitable task models* that fit Linux thread behavior!
- 4. Build and **test** your own in-house tooling.
 - → Or a mess of buggy shell scripts...



LiME in Action

Enter event or section life



Demo: Tracing cyclictest



lime-rtw view /tmp/demo

MPI-SWS

Comma |>> clock_nanosleep (CLOCK_MONOTONIC, absolute) [9999 occurrences] [P] [P] Job Separator By default, cyclictest uses **clock_nanosleep** to enact periodic activations.

|Separator: clock_nanosleep (CLOCK_MONOTONIC, absolute)

Demo: Tracing ping

Can trace arbitrary (and non-RT) tasks.

\$ sudo lime-rtw extract --best-effort -o demo-ping -- ping -c 1000 -i 0.1 contact.mpi-sws.org

PING contact.mpi-sws.org (139.19.171.199) 56(84) bytes of data. 64 bytes from contact.mpi-sws.org (139.19.171.199): icmp_seq=1 ttl=59 time=1.27 ms 64 bytes from swscontact.mpi-sws.org (139.19.171.199): icmp_seq=2 ttl=59 time=1.30 ms 64 bytes from swscontact3.mpi-sws.org (139.19.171.199): icmp_seq=3 ttl=59 time=1.30 ms 64 bytes from swscontact3.mpi-sws.org (139.19.171.199): icmp_seq=4 ttl=59 time=1.32 ms 1.... 64 bytes from swscontact3.mpi-sws.org (139.19.171.199): icmp_seq=1000 ttl=59 time=1.28 ms

--- contact.mpi-sws.org ping statistics ---1000 packets transmitted, 1000 received, 0% packet loss, time 100508ms

rtt min/avg/max/mdev = 1.248/1.309/2.106/0.058 ms

Results saved in demo-ping.

MPI-SWS

lime-rtw view /tmp

-Tasks				1	-Separator	'S
Task ID	Policy	Prio	Comm	Comma	>> <mark>suspe</mark> r	n <mark>sion</mark> [3001 o
>> 494691-0	CFS	-	ping	ping	poll	<mark>g</mark> [999 occur [998 occurren
					1 	
		•			i	Sui
<u>S</u>	uspens	sion S	eparator			
	Each th	iread a	wake-up		- -Separator	• Details
	is a new	w acti	activation.		Separator	: suspension
					Count: 30	001
					 Arrival N	Models:
					 Model #1:	arrival_cur
					Points:	31 (see tab
				I	Model #2:	sporadic
					Minimun	1 Inter-arrıv
					Model #3:	periodic
					Period:	33520000 ns
					Jitter:	118511893 n
					Offset:	42260249008
					I Unt rei	lease

					•	
/כ	d	en	10.	-p	コー	ng
				-		

ccurrences	[P]
rences] [P]	
ices]	

Release Curve

table for complex activation behavior.

rve oles for details)

/al Time: 1100041 ns

ns 338459 ns

- View details | q - Quit

lime-rtw view /tmp/demo-ping

Separator: suspen	nsion Count: 3001	¹ <u>Minimum Separa</u>	ti
Details		Shortest interval in	W
<mark>Summary:</mark> WCET range: 623	86790 - 17705610 ns	→ enables response-t	tir
Arrival Curve Shows the minim	Table: num and maximum len	ng of intervals containing N job re	ele
Number of jobs	Minimum interval	l Maximum interval	
2	1100041	1 98350299	
3	2471339	9 99767058	
	10642605	5 101734014	
4	10042095		
4 5	19561043	3 I 199219047	
4 5 6	10042093 19561043 102438821	3 199219047 1 200622492	
4 5 6 7	10042095 19561043 102438821 107209759	3 199219047 1 200622492 9 202527597	
4 5 6 7 8	19561043 102438821 107209759 116128107	3 199219047 1 200622492 9 202527597 7 299959667	
4 5 6 7 8 9	19561043 102438821 107209759 116128107 202985815	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149	
4 5 6 7 8 9 10	10042093 19561043 102438821 107209759 116128107 202985815 207862863	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149 3 303244402	
4 5 6 7 8 9 10 11	10042093 19561043 102438821 107209759 116128107 202985815 207862863 216781211	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149 3 303244402 1 400678491	
4 5 6 7 8 9 10 11 12	10042093 19561043 102438821 107209759 116128107 202985815 207862863 216781211 303536494	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149 3 303244402 1 400678491 4 402182860	
4 5 6 7 8 9 10 11 12 13	10042093 19561043 102438821 107209759 116128107 202985815 207862863 216781211 303536494 308739186	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149 3 303244402 1 400678491 4 402182860 6 403955856	
4 5 6 7 8 9 10 11 12 13 14	10042093 19561043 102438821 107209759 116128107 202985815 207862863 216781211 303536494 308739186 317657534	3 199219047 1 200622492 9 202527597 7 299959667 5 301397149 3 303244402 1 400678491 4 402182860 6 403955856 4 501379167	

on ("delta min")

hich N jobs arrive. me analysis (RTA)

eases

How Does LiME Work?

MPI-SWS

task = *stream of jobs*

observations.

Key Idea: Job Separators

Target real-time task models: "*a task is a sequence of jobs*"

But Linux is *jobless*: each thread is an *arbitrary* sequence of system calls!

Fundamental Challenge

Where does one job end and the next start?

(1) LiME Insight

On Linux, a job (= one task activation) *always* starts with a *potentially blocking* system call.

A real-time task **must** wait for next event or *passage of time* → *blocking system calls.*

MPI-SWS

(2) LiME Heuristic

Recurrent real-time tasks \approx "doing the same thing over and over"

→ typically every job starts with the same system call, so LiME assumes this.

Job Separator Examples

}

Time-driven: struct timespec next = now(); while (true) { /* periodic activation */ clock_nanosleep(&next); do_something(); timespec_add(&next, PERIOD); }

→ job separator: clock_nanosleep

In both examples:

- one job = one loop iteration
- each job starts with *potentially* blocking system call
- all jobs of one task start with the *same* system call

Event-driven:

int fd = open_udp_socket(PORT);
while (true) {
 /* sporadic activation */
 int nbytes = read(fd, &buf);
 if (nbytes <= 0) break;
 do_something(buf, nbytes);</pre>

→ job separator: read(fd)

Inherent ambiguity:

 task may execute many different system calls in do_something()
 There are potentially many job separator candidates.

Dataflow from Kernel to Model Extraction

LiME design choice: expose job-separator ambiguity to user.

MPI-SWS

user space (Rust)

Supported Task Models

Supported Task Models

Arrival Models

WCET (Liu & Layland, 1973)

WCET(*n*)

(Quinton et al., 2012)

Self-Suspension

MPI-SWS

sporadic

(Mok, 1983)

arrival curves

(Thiele *et al.*, 2000; Richter, 2005)

periodic + jitter & offset

(Liu & Layland, 1973; Leung & Merril, 1980; Audsley *et al.*, 1993)

Execution Time (*measurement-based*)

dynamic (Ming, 1994)

generalized segmented (von der Brüggen et al., 2017)

Why not stick with "tried and true" scalar model parameters?

Example: Linux kernel threads

- migration threads (one per core)
- assist Linux process scheduler
- execute at max. RT priority
- their CPU use must be accounted for

How much CPU use in 300 ms?

WCET	Arrivals	Bound
scalar	scalar	> 1500 ms
curve	scalar	> 340 ms
scalar	curve	≈ 4.3 ms
curve	curve	≈ 0.99 ms

LiME: The Linux Real-Time Task Model Extractor Björn B. Brandenburg* Cédric Courtaud[†] Filip Marković[‡] Bite Ye* *Max Planck Institute for Software Systems, Germany Huawei Technologies, Paris Research Center, France [‡]University of Southampton, United Kingdom

Abstract—We present LIME, a novel dynamic real-time task Abstract—we present LIVIE, a novel dynamic reat-une task model extractor. LIME observes the temporal behavior of Linux mouer extraction. Little observes the temporal behavior of Little real-time threads and automatically maps the observed activity to established real-time task models: sporadic and periodic tasks, upper and lower arrival curves, cumulative execution-time curves, our for a curves, cumulative execution-time curves, and two colf exponentials (dynamic and compared). I the upper and rower arrival curves, cumulative execution-tune curves, and two self-suspension models (dynamic and segmented). LIME and two sen-suspension models (uynamic and segmented). Little runs on unmodified Linux kernels and requires neither knowledge of real-time theory nor familiarity with Linux internals to be used of real-time theory nor taminarity with Linux internats to be used effectively. An extensive evaluation shows LIME to achieve very high inference accuracy—in particular 100% accuracy for common automotive periods—with low kernel overhead, low latency impact, and low processor utilization (at best-effort priority).

I. INTRODUCTION

The recent acceptance of the PREEMPT_RT patch into the mainline Linux kernel, after 20 years of development [58], has made official what has long been true: Linux is a major platform for hosting modern, complex real-time workloads across various industries. Notable examples feature demanding application domains such as automotive systems [52], autonomous vehicles [30, 31], unmanned aerial vehicles (UAVs) [19, 27, 35], and spacecraft [37], in particular crewed rockets [59], NASA's

Mars helicopter *Ingenuity* [57], and tens of thousands of Linux systems deployed in orbit as part of *Starlink* constellations [54]. As noted recently by Erik Vallow, a representative of the RTOS vendor LYNX Software Technologies, in a retrospective on the evolving role of Linux in the aerospace and defense industries [56]: "Linux has become a formidable contender in safety-critical systems due to advancements in real-time of meeting the demands of real-time applications on its own, is increasingly a factor favoring the consolidation of AI-enabled or otherwise GPU-accelerated real-time workloads on Linux. However, while the popularity of Linux as a versatile and feature-rich RTOS has soared in the real-time systems industry,

there is a growing disconnect with the analytical foundations studied in the scientific literature on real-time systems. Rooted in abstract system models and high-level mathematical descriptions of workloads, state-of-the-art methods for establishing temporal guarantees are far removed from the engineering realities of a low-level embedded Linux environment.

It stands to reason, then, that only a diminishingly small fraction of the many real-time workloads deployed on Linux

over the past decade have been modeled and formally evaluated using published schedulability analyses. A major contributor to this disconnect is the lack of tool support: without automated

system introspection tools, engineers interested in formally characterizing the timing properties of a real-time workload running on Linux require dual expertise in both Linux implementation details, particularly the kernel's low-level tracing facilities and system-call interface, and state-of-the-art temporal modeling and analysis techniques. This, along with the associated *manual* effort that would be required (and not just once, but repeatedly as systems evolve to meet changing requirements), presents a formidable barrier to the widespread adoption of state-of-the-art real-time analysis techniques in a Linux context. Could the schedulability analysis of Linux workloads be

automated and thus made easily accessible to non-experts? Motivated by this question, we address the first, fundamental problem that precedes any practical analysis: the *automated* modeling of real-time tasks deployed on *unmodified* Linux kernels. While applications developed using higher-level model-first approaches [9, 24, 44], or using programming languages with explicit timing semantics [7, 10, 11, 15, 26, 28, 42, 43], can certainly be *compiled down* to Linux binaries and executed efficiently (*i.e.*, model-driven engineering), the converse is far from obvious: Is it possible to extract high-level temporal *analysis simply by observing their low-level runtime behavior? capabilities and reliability.* [...] *Linux is increasingly capable of meeting the demands of real-time applications on its and capabilities and reliability.* [...] *Linux is increasingly capable is it possible to do this for black-box threads (i.e., without access to source code)?* Fully *automatically, without with* of meeting the aemands of real-time applications on its own, access to source code)? Fully automatically, without user reducing the need for a separate RTOS in some cases." In guidance, annotations, or specifications of intended timing? And can it be done in situ on a target embedded platform. *reaucing the need for a separate RTOS in some cases.*" In addition, the availability of drivers for high-performance GPUs is increasingly a factor favoring the consolidation of AL applied without unduly perturbing the timing of the real time theory and without unduly perturbing the timing of the real time theory and addition. without unduly perturbing the timing of the real-time threads? This paper. We show that the answer to each of these questions

is 'yes' by presenting LIME, the Linux real-time task Model Extractor (Sec. III). LIME is a dynamic introspection tool

black-box workload (in Fig. 1, a thread implementing periodic activations with clock_nanosleep), LIME uses eBPF to observe key scheduling events and system calls throughout the

Evaluation in the Paper

that maps sequences of low-level thread-kernel interactions (Sec. II-A), observed via the kernel's *eBPF* tracing facility, to task models from the real-time scheduling literature (Sec. II-B). Fig. 1 illustrates the entire pipeline: Given an arbitrary *^{†‡} Authors are listed in alphabetical order, ^{†‡} This work was carried out while affiliated with the Max Planck Institute for Software Systems, Germany. Whole system (Inset 1). After reconstructing the timeline for based on its builtin understanding of Linux system call *^{†‡} Authors are listed in alphabetical order. ^{†‡} This work was carried out while affiliated with the Max Planck Institute for Software Systems, Germany. based on its builtin understanding of Linux system call

LiME is Accurate and Effective

The ROSACE Case Study: From Simulink Specification to Multi/Many-Core Execution

Claire Pagetti*, David Saussié[†], Romain Gratia*, Eric Noulard*, Pierre Siron* [†] Polytechnique Montréal - Canada *ONERA - Toulouse, France

Abstract—This paper presents a complete case study - named **ROSACE** for Research Open-Source Avionics and Control Engineering - that goes from a baseline flight controller, developed in MATLAB/SIMULINK, to a multi-periodic controller executing on a multi/many-core target. The interactions between control steps, in particular by investigating several multi-periodic configurations. We deduced ways to improve the discussion between engineers in order to ease the integration on the target. The

- **P3** : rise time, that is the time it takes to rise from 10% to 90% of the steady-state value;
- **P4** : steady-state error, that is the difference between the input and the output for a prescribed test input as $t \to \infty$.

and computer engineers are highlighted during the development The time-domain performance properties are illustrated in the figure 1 for a step input. At this stage, these properties are analyzed through SIMULINK simulations.

Time

Pagetti et al., RTAS 2014 Settling time

MPI-SWS

- <u>100% Inference Accuracy</u>
- \rightarrow automotive periods
- \rightarrow random periods (ms-granularity)

<u>ROSACE Case Study</u>

- → *longterm* observation
- → detect subtle timing drift
- \rightarrow validate fix

LiME Causes Only Low Overhead

<u>CPU overhead</u> Kernel (eBPF): < 2.5% (of 4 cores) User space: < 2.0% (*on a Raspberry Pi* 4 "Model B")

cyclictest latency benchmark → avg. latency impact $\approx 25\mu s$

Redis key-value server
→ minor throughput impact

Throughput (ops/s) 0 05 10

Conclusion

What Can LiME Do For You?

If you are doing systems research:
→ *Integrate schedulability analysis (almost) "for free"*

If you are working with a **real system**: → *Monitor, understand, debug, validate its timing behavior*

If you work on schedulability analysis:
→ Get some real models rather than just random numbers
→ Demonstrate that your model assumptions are viable

If you **teach**: → *Let students explore RT foundations hands-on*

MPI-SWS

blackbox threads running on Linux

> → highly accurate → fully automated \rightarrow low overheads → online inference

Thank you for your attention! Any Questions?

MAX PLANCK INSTITUTE FOR SOFTWARE SYSTEMS

established realtime task models

European Research Council Established by the European Commission