



Virtual Machine Boot-up Analysis

Abderrahmane Benbachir

Dec 12, 2016

École Polytechnique de Montréal
Laboratoire **DORSAL**

Agenda

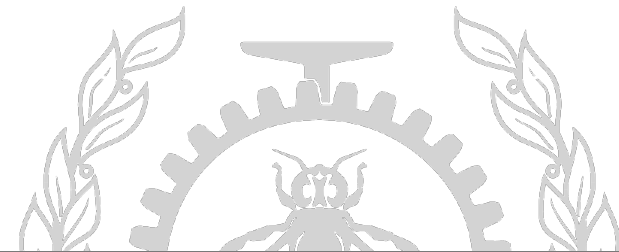
Introduction

Research objectives

Investigation

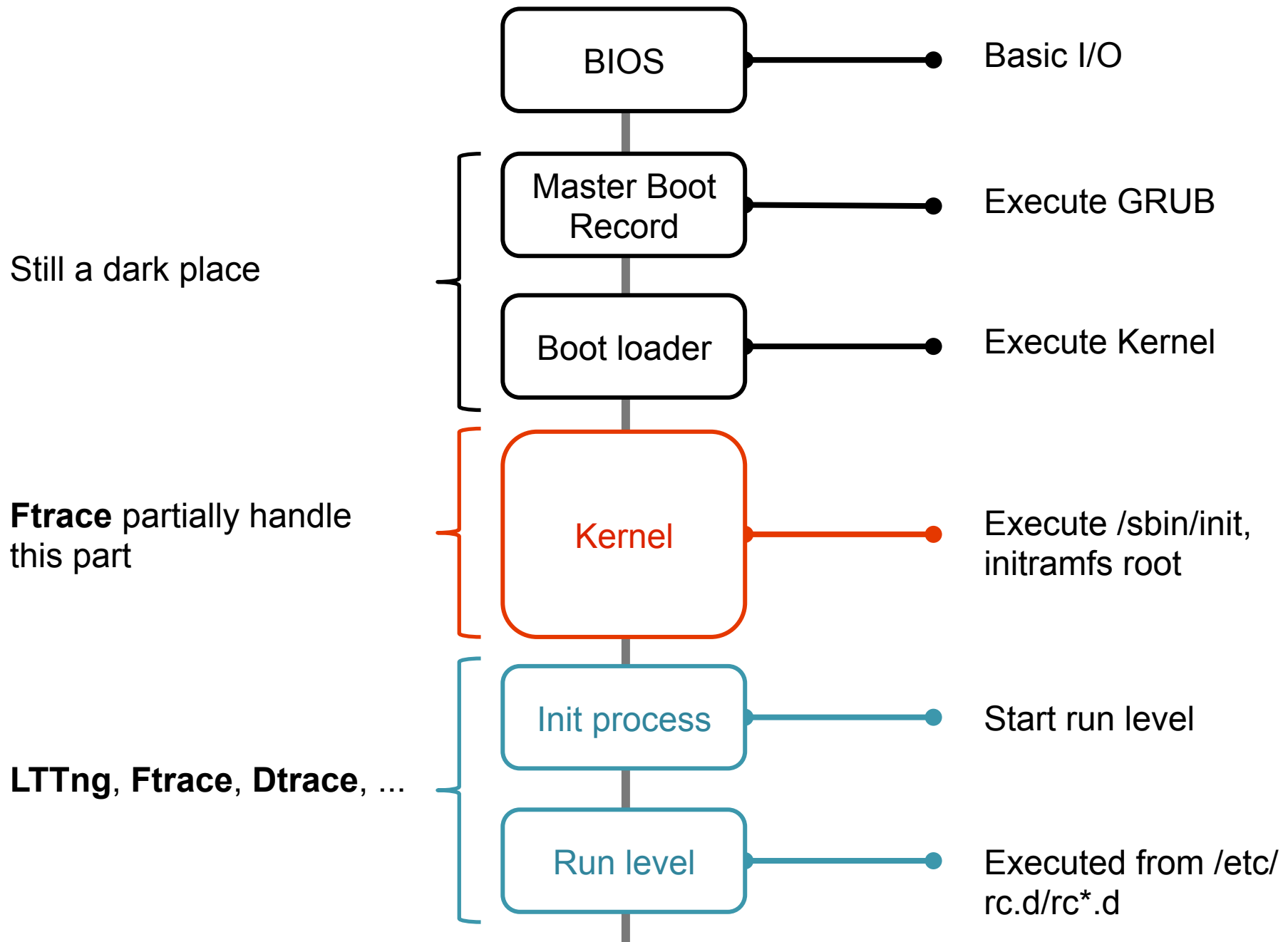
Results

Future work



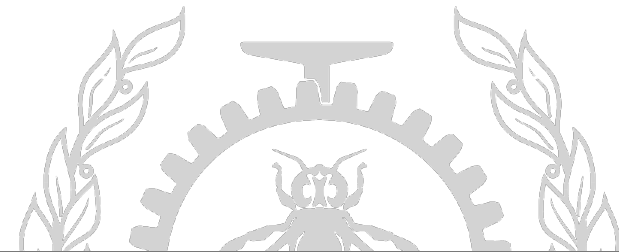
Introduction

Boot-up process



Research objectives

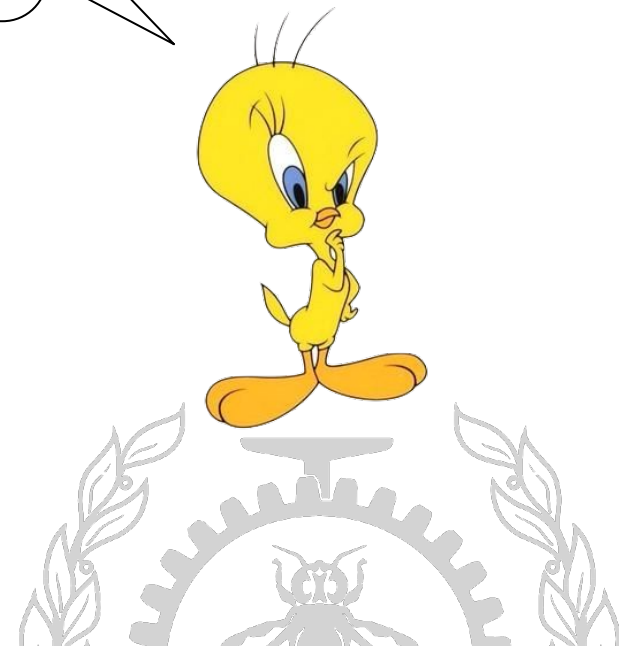
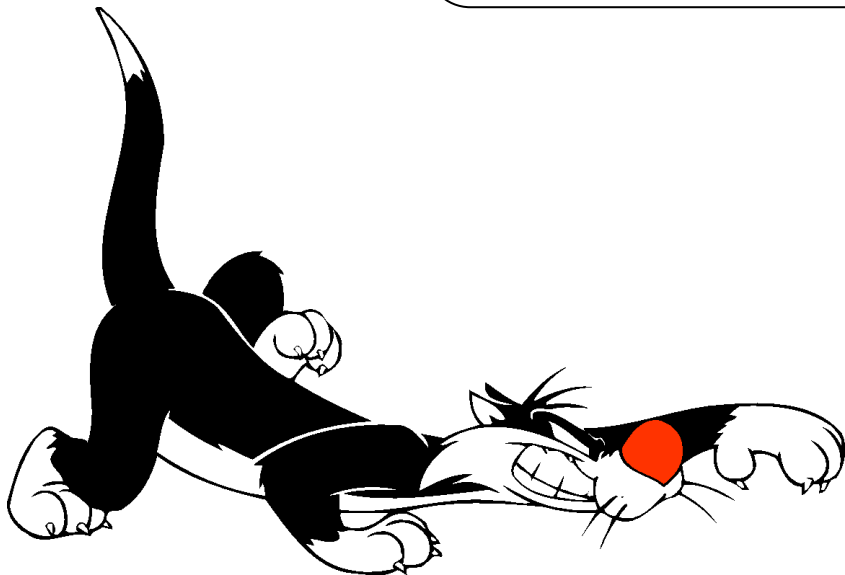
- Analyze virtual machines early boot-up
- Detect boot up issues



Investigation

We already know how to trace: **LTng**, **Ftrace**, **Dtrace**, **Perf**.

Where & How to store traces during bootup ?



Source : <https://s-media-cache-ak0.pinimg.com/736x/99/66/28/996628fa4292de7175d7893ca31d7796.jpg>

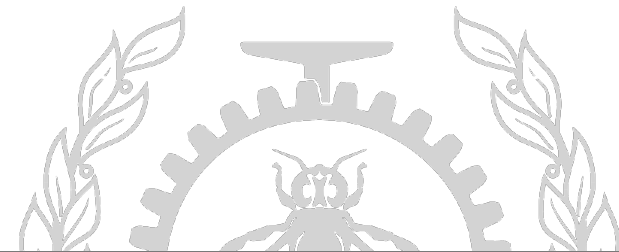
Keep them inside **Guest**

- Recording Timestamp
- Synchronization
- Buffer size
- Memory or Disk

vs

Offload them to the **Host**

- Guest or Host Timestamp ?
- Synchronization
- Transfer Channels (Network, Shared Memory, Paravirt)
- Transfer Speed \$\$\$



Investigation

Where & How to store traces during boot up ?

Keep them inside
Guest

vs

Offload them to the
Host

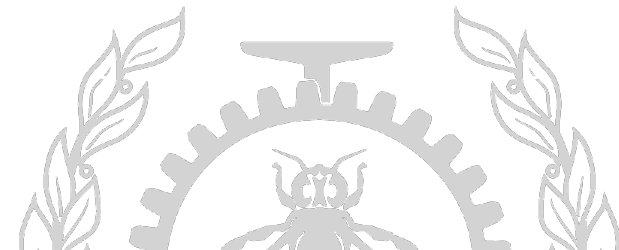
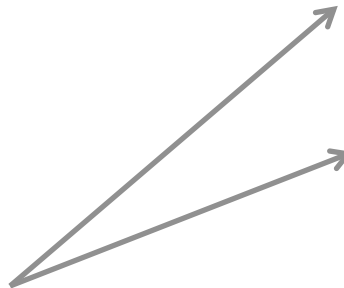
Light weight Tracers

Network : Sockets

Shared Memory: Qemu Hypertrace

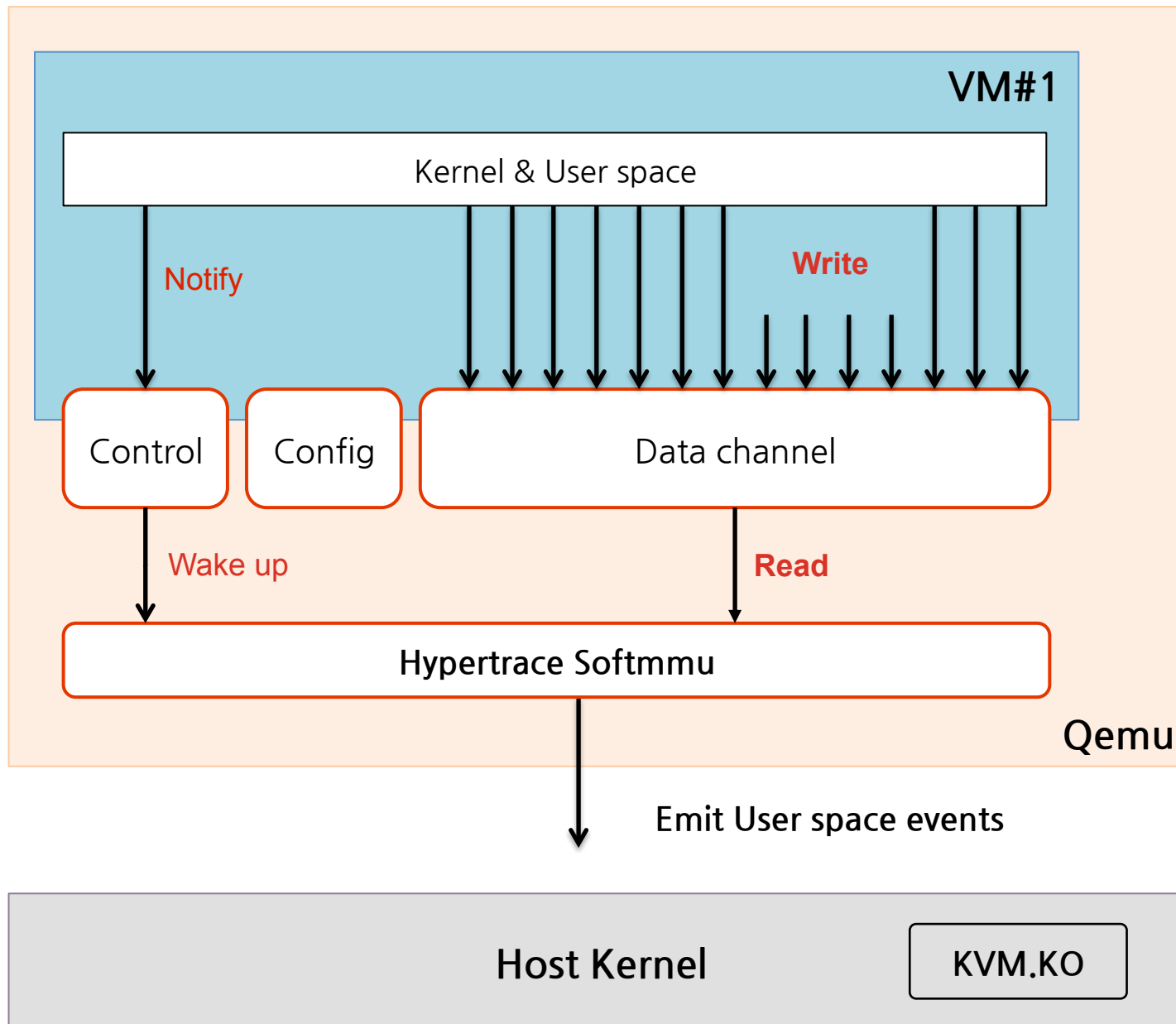
Paravirtualization API: Hypercalls

Must be explored

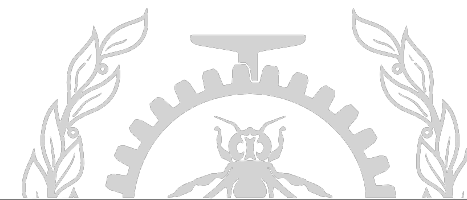


Qemu Hypertrace

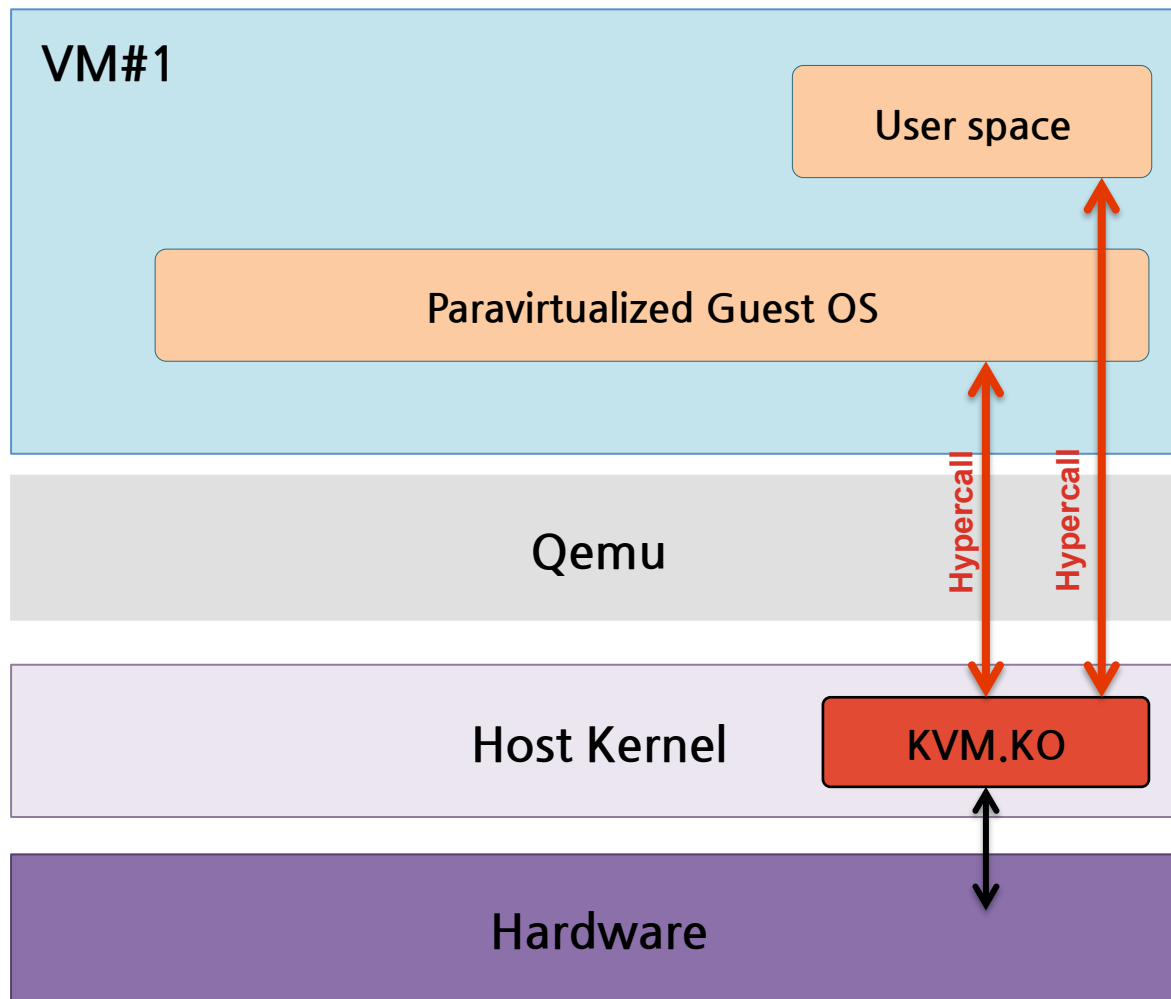
By Lluís Vilanova



**Shared Memory
Btw
Guest & Qemu**

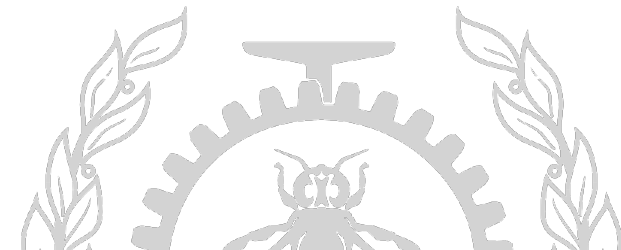


Hypercall



**Direct Execution
of
Guest Requests**

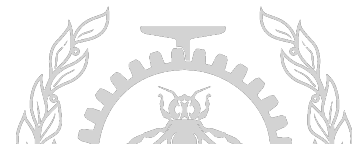
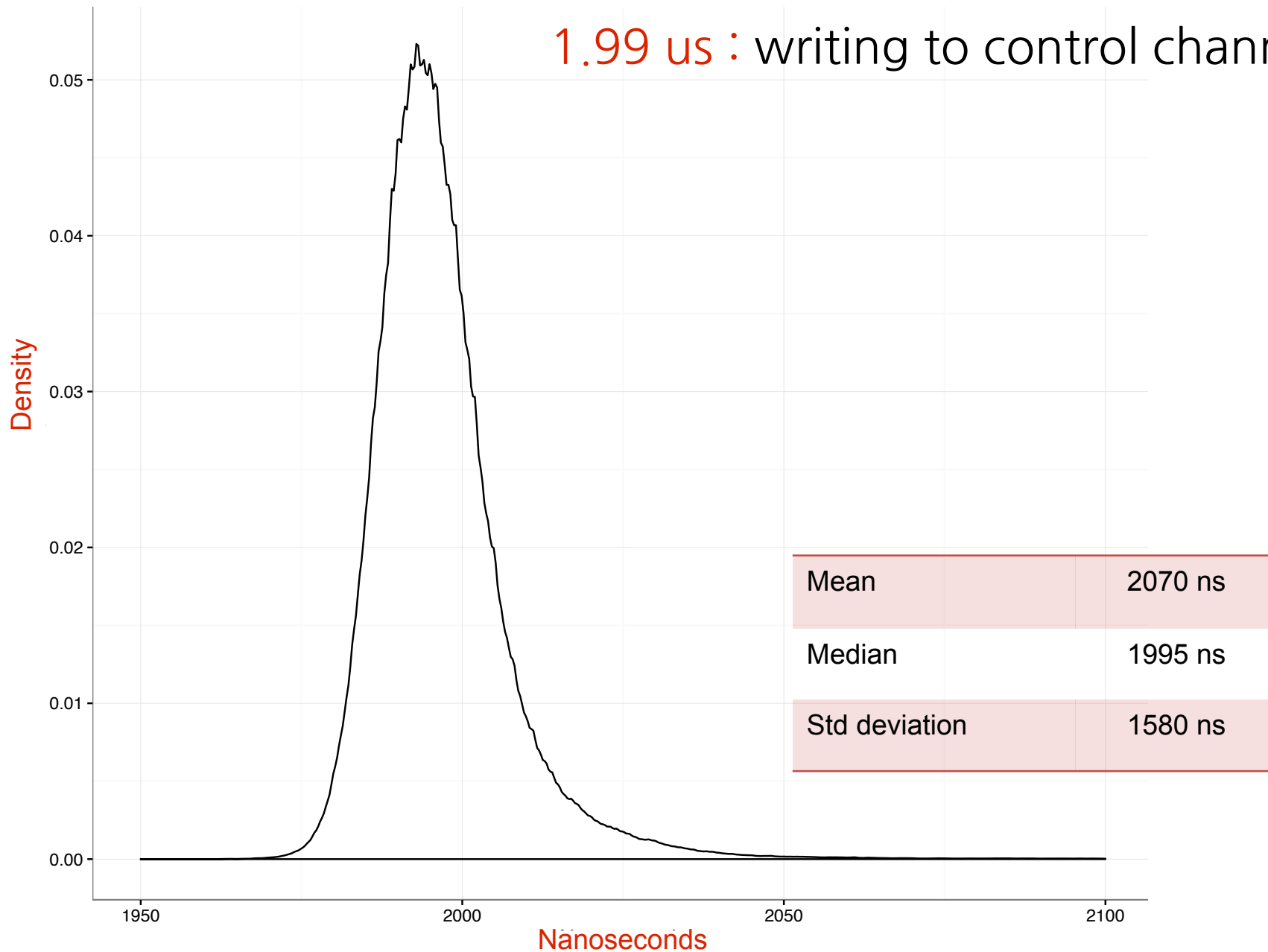
This layer is not involved



Benchmarks

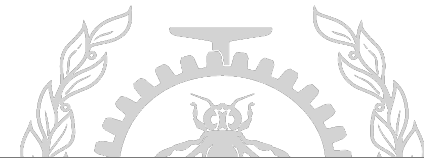
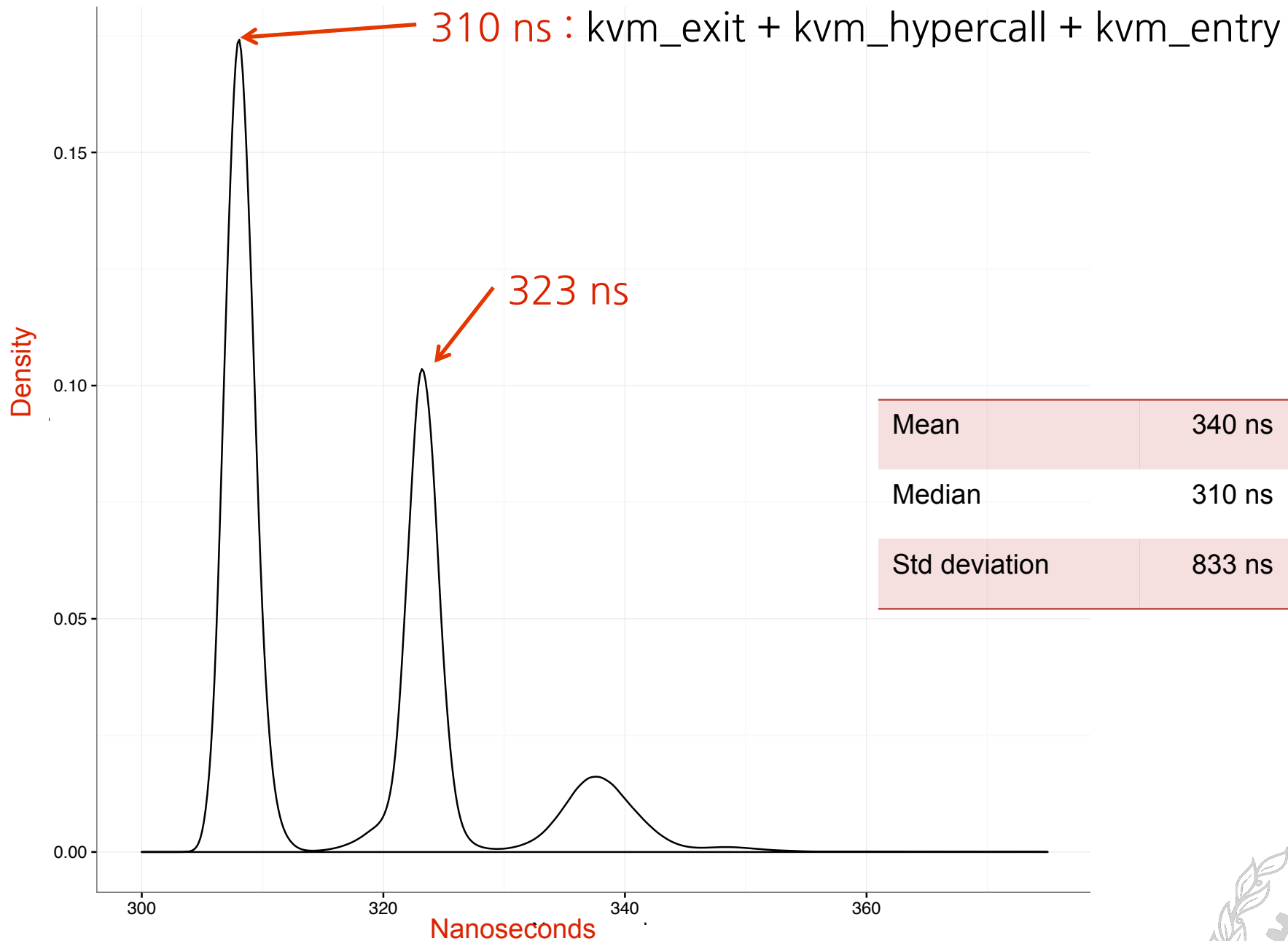
Qemu Hypertrace Overhead (KVM Enabled)

1.99 us : writing to control channel



Benchmarks

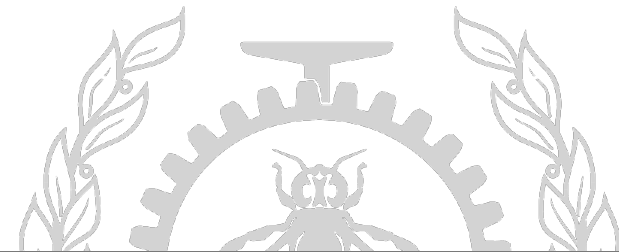
Hypercall Overhead



Hypercalls Implementation

How to trace through Hypercalls

- Hook to Ftrace “function graph” callback for entry & exit
- Only trace the Host
- Use only host Timestamp : No synchronization required
- Dump kernel symbols to map each function names



Hypercalls Implementation

hypercall payload = 5 args x 64 bits = 40 Bytes

hypercall(nr, a0, a1, a2, a3)

Hypercall number

Entry

[13:59:00.035571626] (+0.000000494) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579284592, a1 = 0, a2 = 0, a3 = 0 }

[13:59:00.035572670] (+0.000000519) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579288208, a1 = 1, a2 = 4198, a3 = 3 }

[13:59:00.035573178] (+0.000000508) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579292656, a1 = 1, a2 = 5227, a3 = 2 }

[13:59:00.035573686] (+0.000000508) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579922416, a1 = 1, a2 = 6251, a3 = 1 }

[13:59:39.016994097] (+0.000000482) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4197127, a1 = 0, a2 = 0, a3 = 0 }

[13:59:00.035584321] (+0.000000531) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579265552, a1 = 0, a2 = 0, a3 = 0 }

[13:59:00.035586871] (+0.000000490) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579284528, a1 = 0, a2 = 0, a3 = 0 }

[13:59:00.035587383] (+0.000000525) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579284592, a1 = 1, a2 = 40, a3 = 4 }

[13:59:39.016994570] (+0.000000473) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4197127, a1 = 1, a2 = 0, a3 = 0 }

[13:59:39.016995047] (+0.000000477) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4196150, a1 = 0, a2 = 0, a3 = 0 }

[13:59:39.016996002] (+0.000000481) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4196518, a1 = 0, a2 = 0, a3 = 0 }

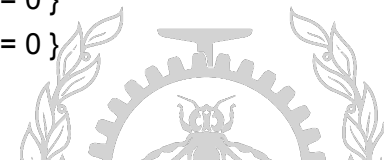
[13:59:00.035585361] (+0.000000503) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071579284592, a1 = 0, a2 = 0, a3 = 0 }

[13:59:00.035585863] (+0.000000502) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 1000, a0 = 18446744071587317584, a1 = 0, a2 = 0, a3 = 0 }

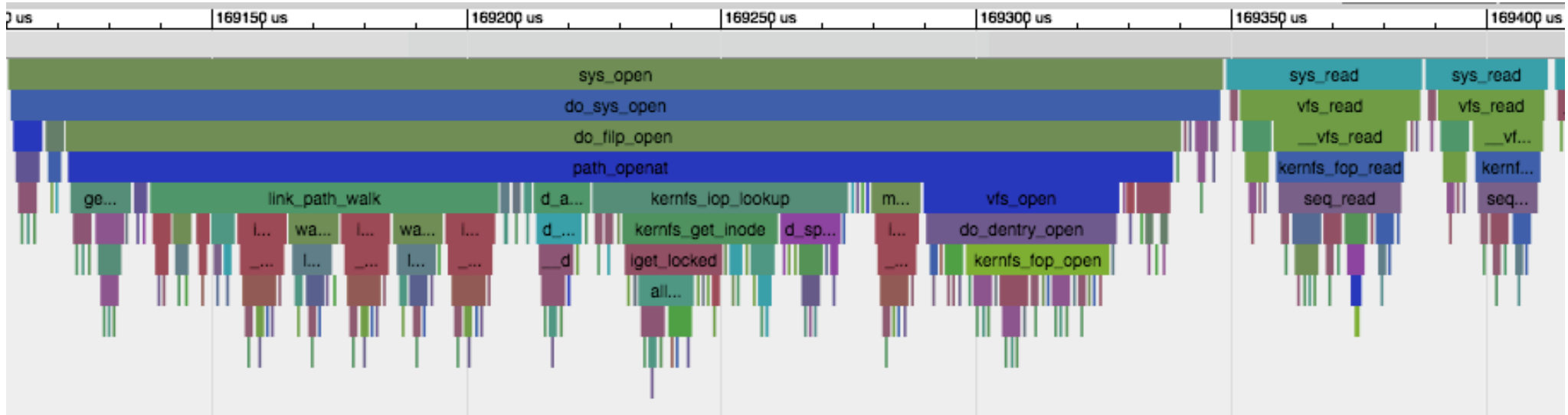
[13:59:39.017128748] (+0.000000734) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4196518, a1 = 1, a2 = 0, a3 = 0 }

[13:59:39.017129250] (+0.000000502) abder-pc kvm_x86_hypercall: { cpu_id = 3 }, { nr = 2000, a0 = 4197127, a1 = 1, a2 = 0, a3 = 0 }

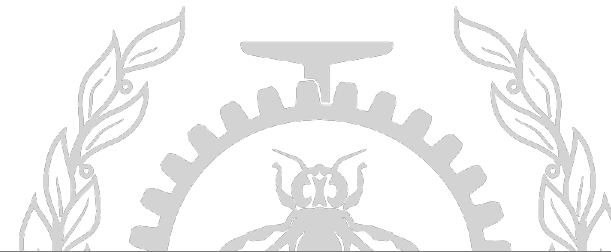
Exit



Callstacks of Guest Kernel space

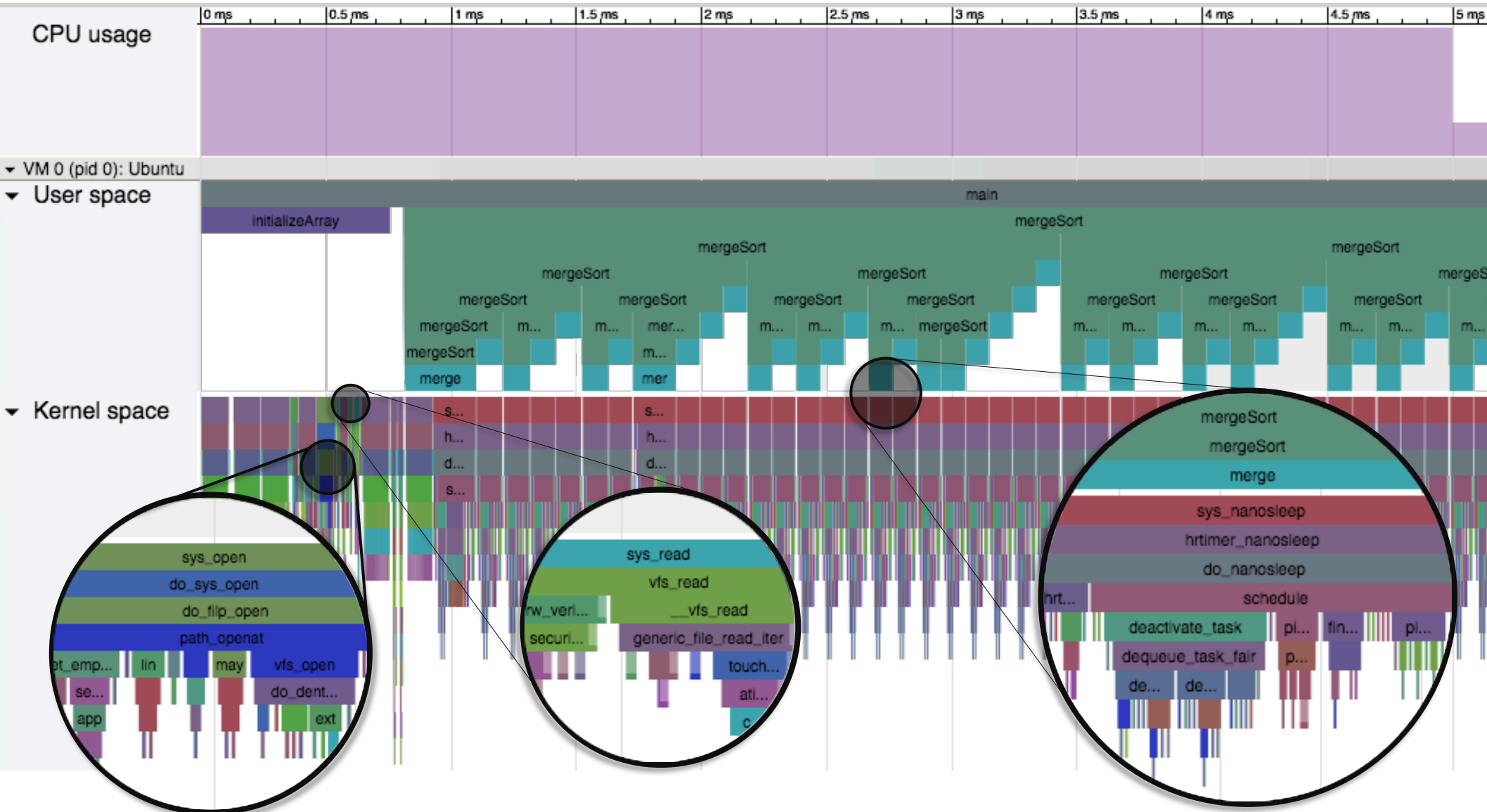


← 300 us →

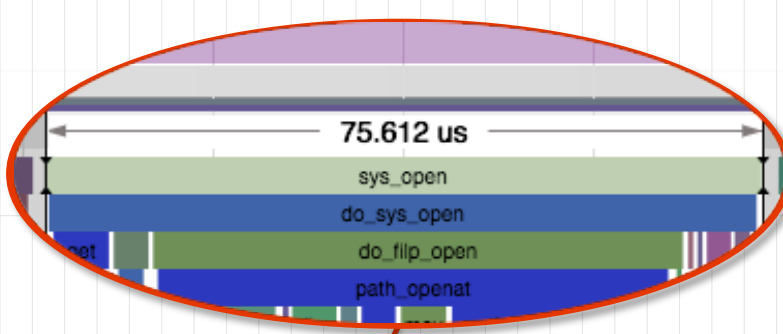
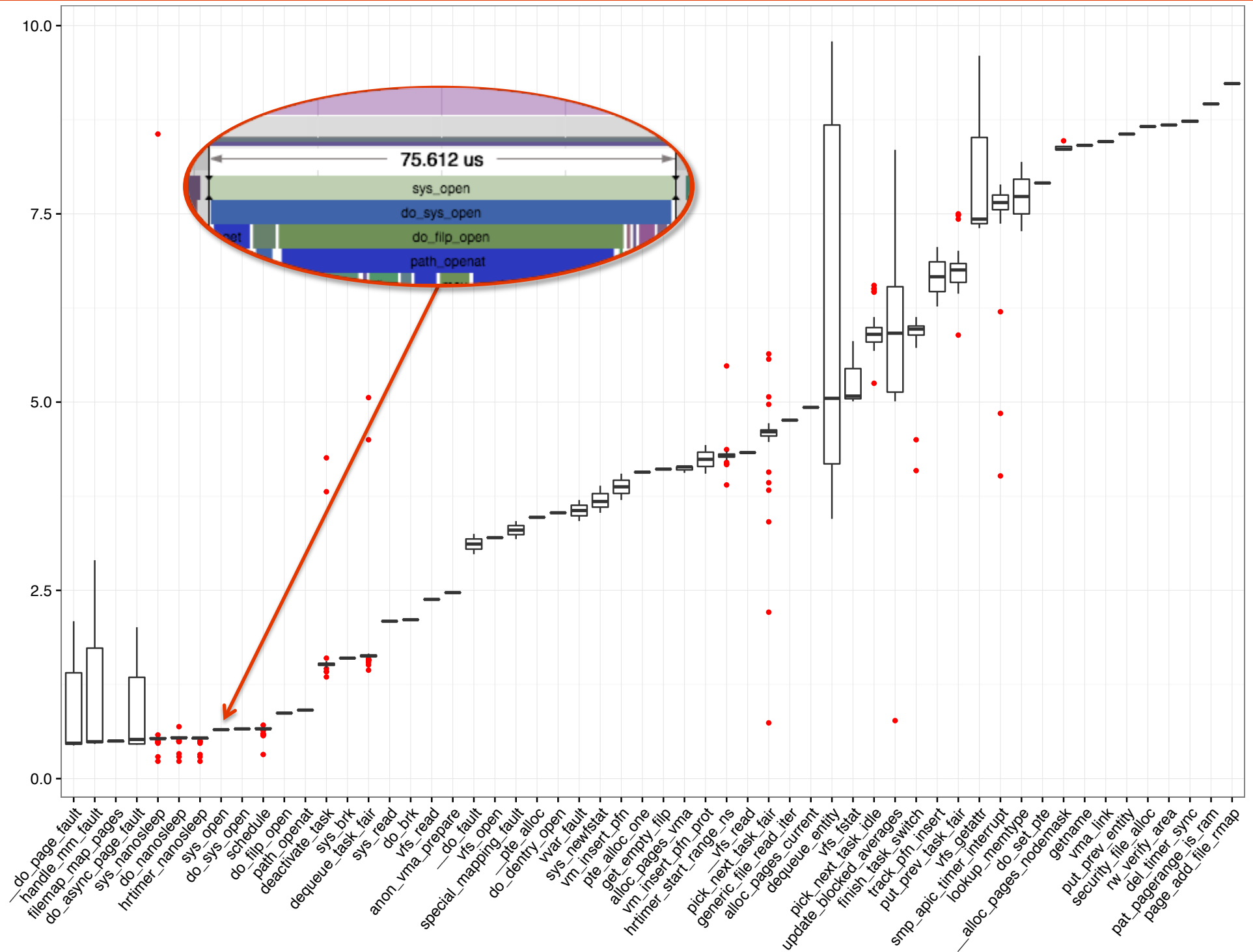


Hybrid tracing technique

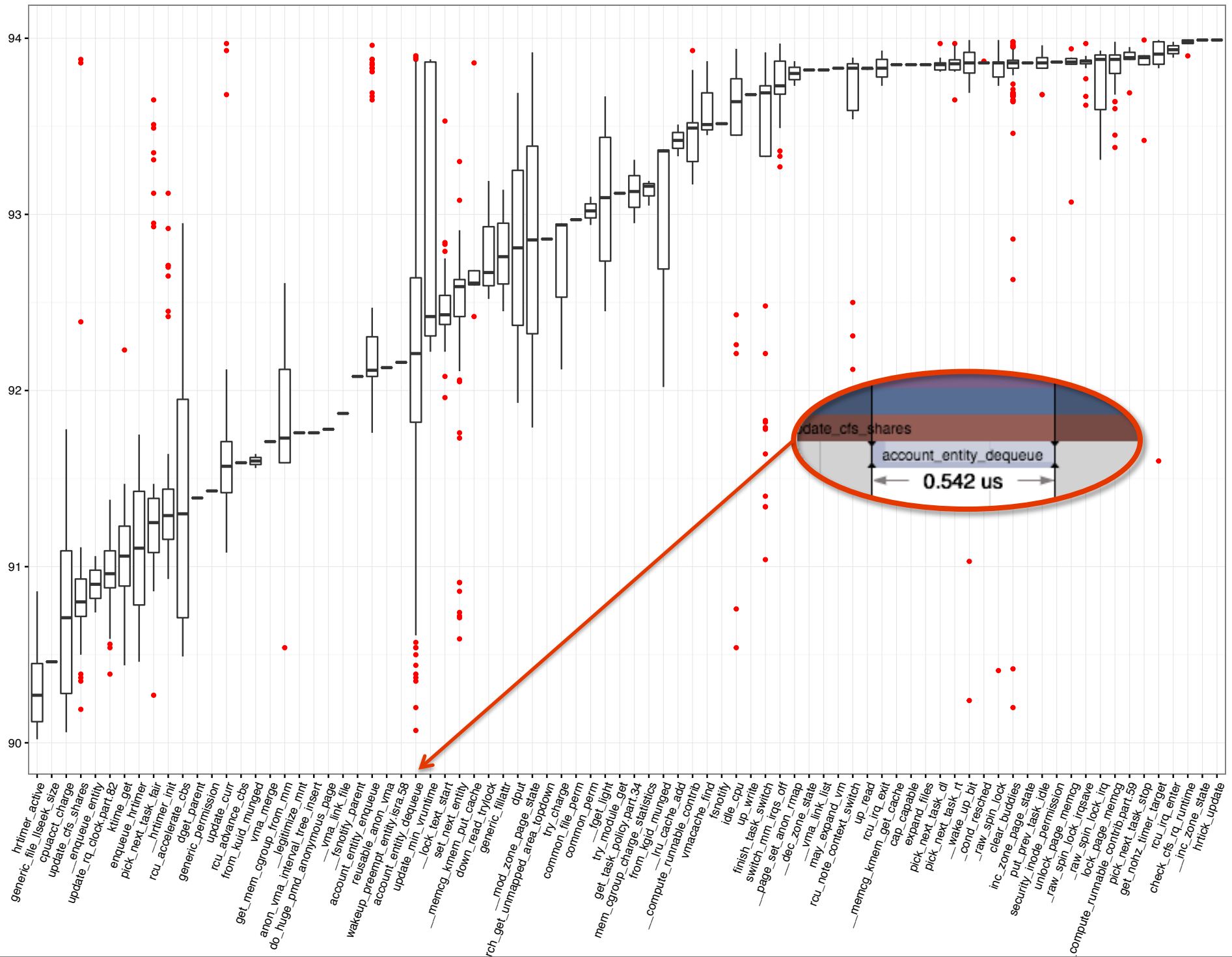
Kernel space + User space



Overhead %



Overhead %



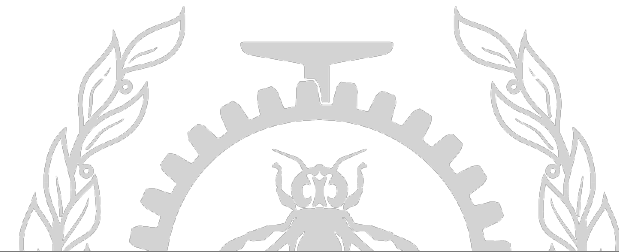
update_cfs_shares

account_entity_dequeue

0.542 us

Future Work

- Apply hypercall technique to trace full kernel boot up
- Explore this solution with nested virtual machine
- Can we trace the Boot loader too ?

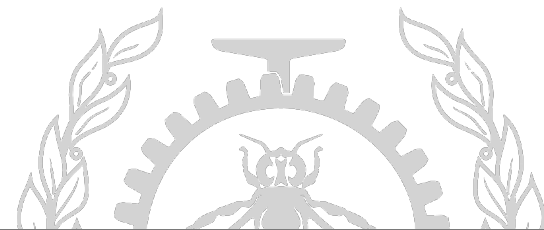


Questions

Abderrahmane.benbachir@polymtl.ca

<https://github.com/abenbachir>

Thanks to Professor Michel Dagenais and our partners
EfficiOS and Ericsson.



References

Hypercall Implementation : <https://gist.github.com/abenbachir/344822b5ba9fc5ac384cdec3f087e018>

QEMU Hypertrace Patches: <http://patchwork.ozlabs.org/project/qemu-devel/list/?state=&q=Hypertrace&archive>

