



Mitigating Excessive Pause-Loop-Exiting in VM-Agnostic KVM

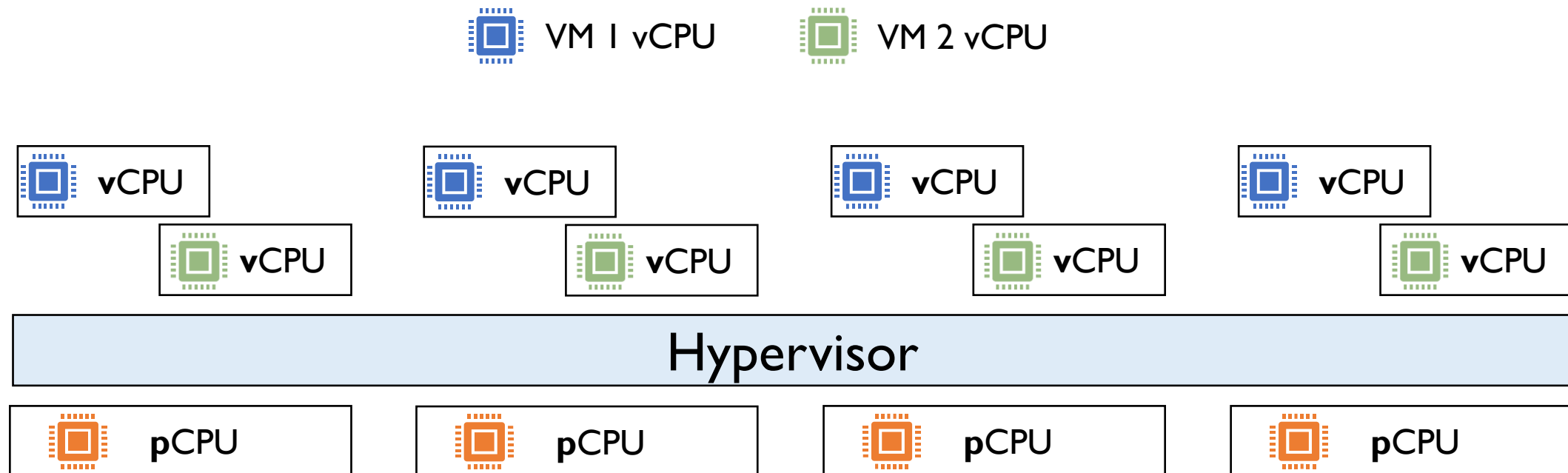
Kenta Ishiguro / Keio University

About Me

- 3rd year PhD student at Keio University, Japan
- Research interests:
 - performance and security of hypervisors
 - particularly KVM

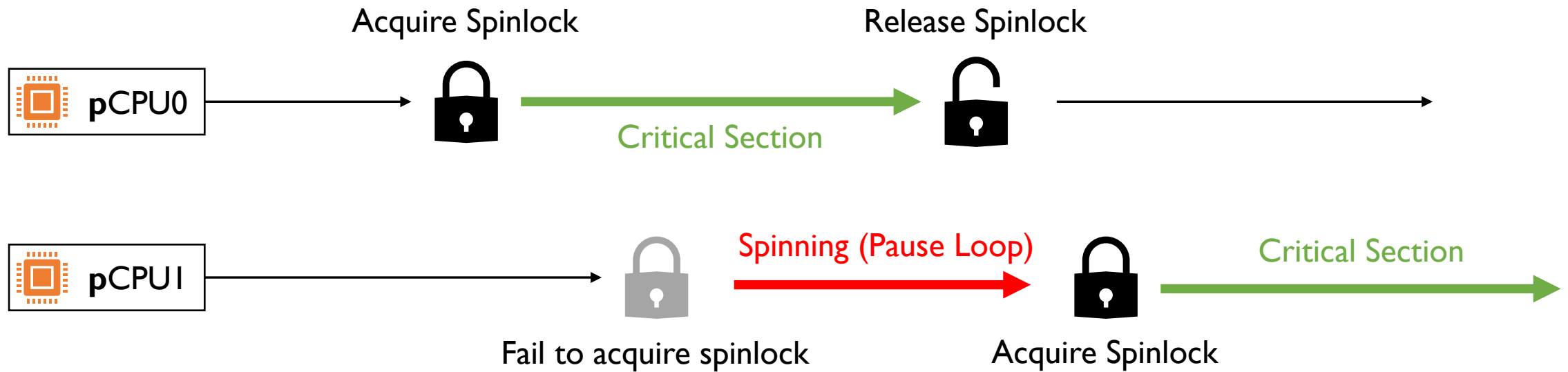
Oversubscribing **virtual** CPUs

- Enables better hardware utilization
- Requires multiplexing of virtual CPUs (vCPUs) on physical CPUs (pCPUs)



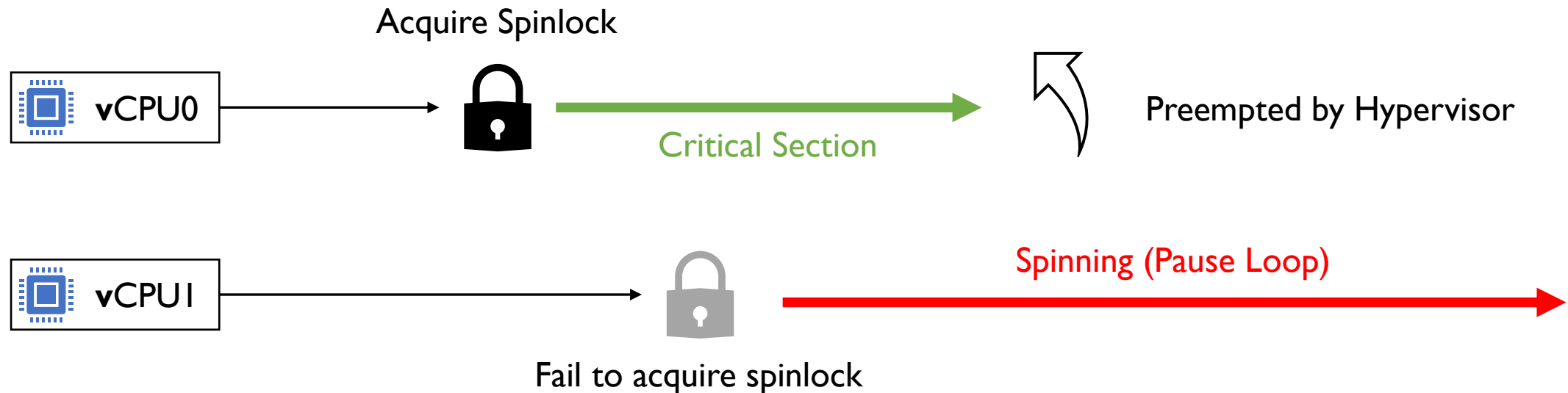
Critical section on pCPU

- Operating System assumption:
 - pCPUs are **always active**



Excessive vCPU spinning

- Oversubscription incurs **excessive spinning**
- **vCPU pre-emption** by Hypervisor **violates OS assumption**
- Lock-holder pre-emption problem



Improving guest performance by solving excessive vCPU spinning

- Problem
 - vCPU **wastes** its execution time by **excessive spinning**
- Ideal case
 - Hypervisor knows which vCPU should be scheduled right now
- But, hard to solve excessive vCPU spinning due to **semantic gap**
 - between KVM and Linux scheduler
 - between KVM and guest VMs

Improving guest performance by solving excessive vCPU spinning

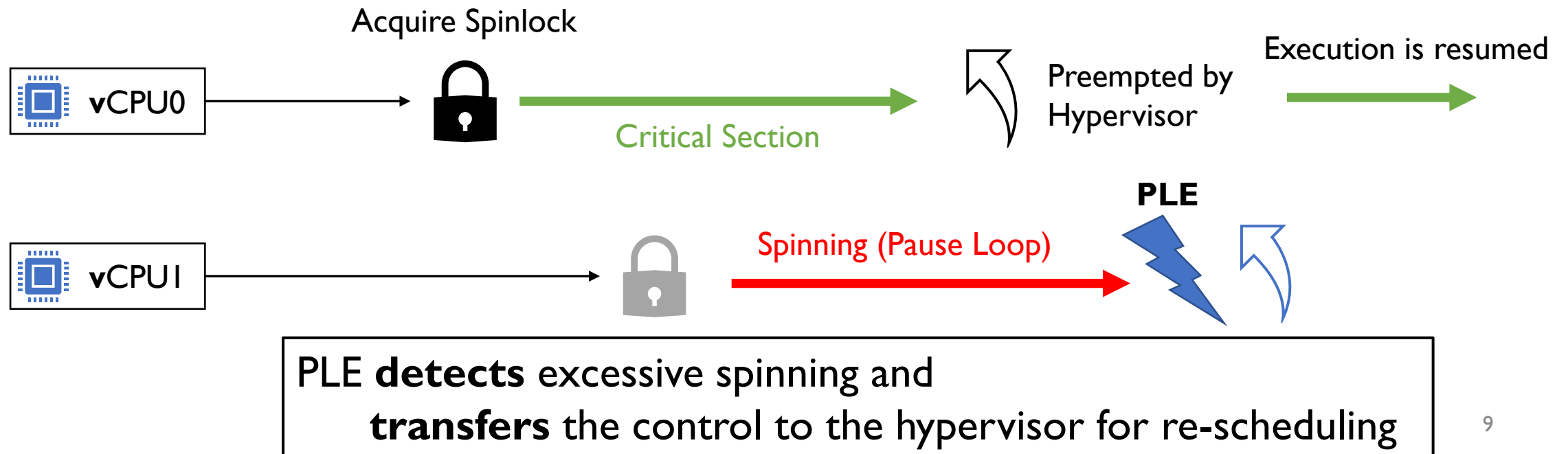
- Problem
 - vCPU **wastes** its execution time by **excessive spinning**
- Ideal case
 - Hypervisor
 - Boosting a target vCPU can be impeded by the Linux scheduler to keep fairness between vCPUs
- But, hard to solve excessive vCPU spinning due to **semantic gap**
 - between KVM and Linux scheduler
 - between KVM and guest VMs

Improving guest performance by solving excessive vCPU spinning

- Problem
 - vCPU **wastes** its execution time by **excessive spinning**
- Ideal candidate
 - Hypervisor for boosting
 - Hard to build comprehensive vCPU candidate set
- But, hard to solve excessive vCPU spinning due to **semantic gap**
 - between KVM and Linux scheduler
 - between KVM and guest VMs

KVM approach with hardware feature

- KVM leverages Pause Loop Exiting (**PLE**) on Intel x86
 - **detects** excessive vCPU spinning
 - gives hypervisors a chance to **re-schedule** vCPUs



KVM's strategy to suppress PLE events

- Reschedule PLE-ed vCPU to another preempted vCPU
 - in `kvm_vcpu_on_spin`
- Co-operative rescheduling with Linux scheduler
 - leverages `yield_to` provided by CFS scheduler
 - makes a request to Linux scheduler to yield and boost vCPUs
- Selects a vCPU to boost in round-robin from candidates
 - resource-waiter, lock-waiter, and IPI-receiver vCPUs are candidates

What happens in the worst case?

KVM trace with running two 8-vCPU VMs simultaneously

```
CPU 2/KVM 2679 [006] 276.731784: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731819: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731823: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731845: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731849: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731871: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731875: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731897: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731901: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731923: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731927: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731949: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731954: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731975: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731980: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732001: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732006: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732027: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732032: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732053: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732058: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
```

PLE events occur continuously (> 100 times)

What happens in the worst case?

KVM trace with running two 8-vCPU VMs simultaneously

```
CPU 2/KVM 2679 [006] 276.731784 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731819 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731823 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731845 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731849 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731871 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731875 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731897 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731901 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731923 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731927 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731949 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731954 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731975 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731980 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732001 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732006 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732027 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732032 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732053 kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732058 kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
```

Continuous PLE occur in the short period

What happens in the worst case?

KVM trace with running two 8-vCPU VMs simultaneously

```
CPU 2/KVM 2679 [006] 276.731784: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731819: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731823: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731845: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731849: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731871: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731875: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731897: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731901: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731923: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731927: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731949: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731954: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731975: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731980: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732001: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732006: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732027: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732032: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732053: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732058: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
```

At the same code location

What happens in the worst case?

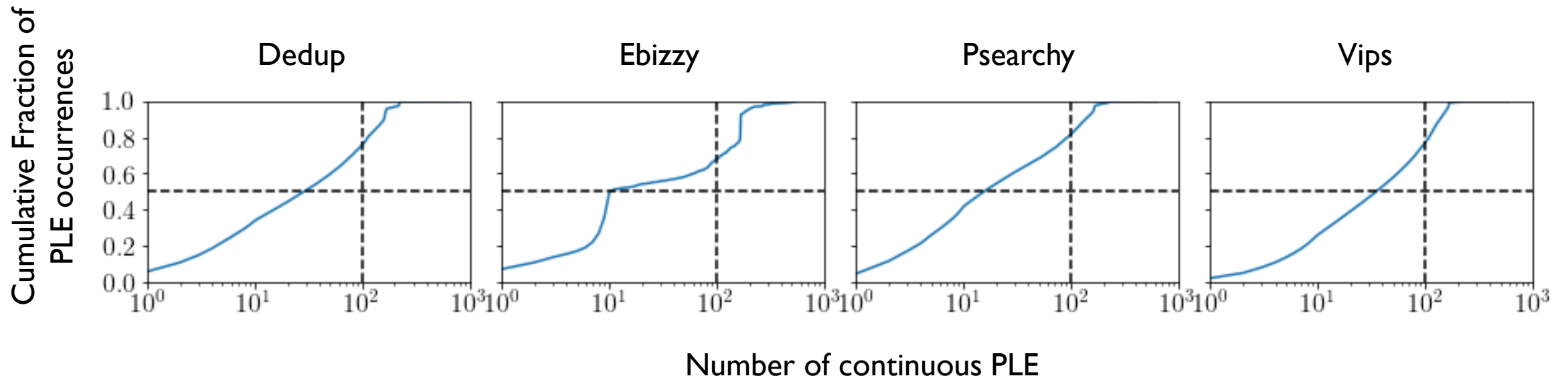
KVM trace with running two 8-vCPU VMs simultaneously

```
CPU 2/KVM 2679 [006] 276.731784: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731819: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731823: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731845: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731849: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731871: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731875: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731897: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731901: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731905: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731909: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731913: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731917: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731921: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731925: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731929: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731933: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731937: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731941: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731945: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731949: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731953: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731957: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731961: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731965: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731969: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731973: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731977: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731981: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731985: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731989: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.731993: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.731997: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
CPU 2/KVM 2679 [006] 276.732001: kvm:kvm_entry: vcpu 2
CPU 2/KVM 2679 [006] 276.732005: kvm:kvm_exit: vcpu 2 reason PAUSE_INSTRUCTION rip 0xffffffffb72e29e3 info 0 0
```

- PLE events occur continuously (> 100 times)
 - in the short period
 - at the same code location

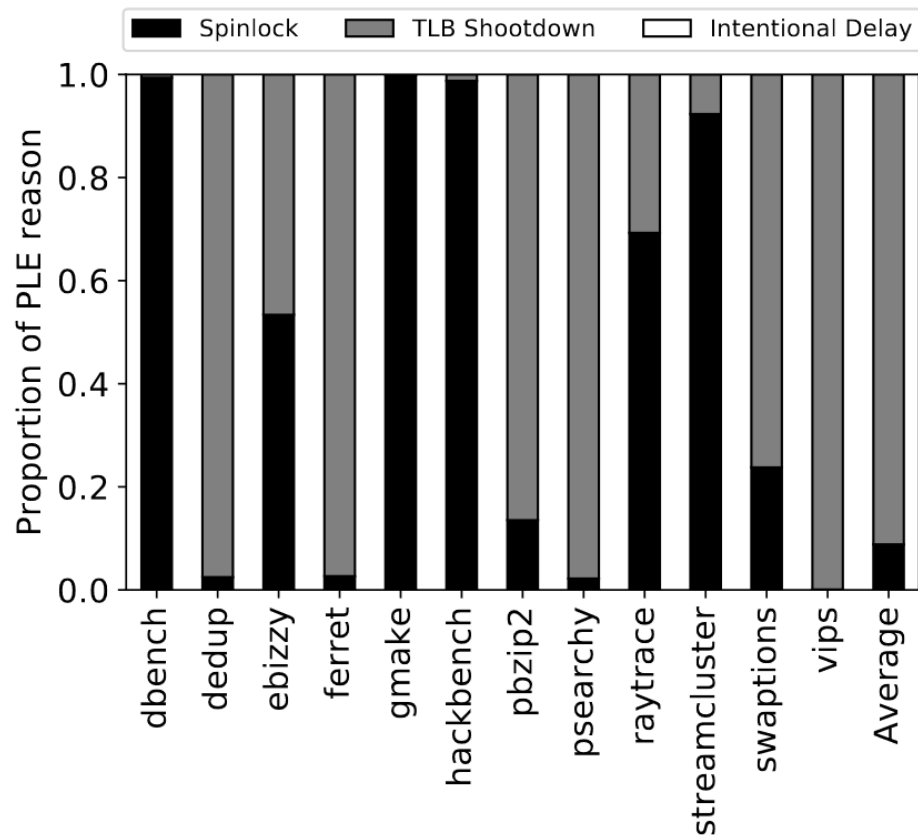
Continuous PLE events are **NOT** rare

- Most number of PLE events are from continuous PLE events
- Cause of PLE events: both spinlock and TLB shutdown



Continuous PLE events are **NOT** rare

- Most number of PLE events are from continuous PLE events
- Cause of PLE events: both spinlock and TLB shutdown



Functionality	Function causing PLE events
spinlock	native_queued_spin_lock_slowpath queued_write_lock_slowpath queued_read_lock_slowpath try_to_wake_up etc...
TLB shutdown	smp_call_function_many smp_call_function_single

Why PLE events occur continuously?

- Lost opportunity and Overboost
 - Comprehensive approach to identify the root cause of PLE does not exist
 - Mitigation: **Strict Boost**
- Scheduler mismatch
 - Linux scheduler **ignores hints** from hypervisor
 - Mitigation: **Debooster**

How to select candidate vCPU

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**
- Boost halted vCPUs where they have received an IPI [2019]
 - Optimization against **TLB shutdown**

How to select candidate vCPU

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**
- Boost preempted vCPUs where they have received an IPI [2019]

Insights

on against **TLB shutdown**

- Spinlock still incurs continuous PLE events
- vCPUs in user mode also need to ack when TLB shutdown happens
- No need to boost IPI-receivers when lock-holder preemption happens

How to select candidate vCPU

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**

- Boost preempted vCPUs again

Insights

due to scheduler mismatch (described later)

- **Spinlock** still incurs continuous PLE events
- vCPUs in user mode also need to ack when TLB shutdown happens
- No need to boost IPI-receivers when lock-holder preemption happens

How to select candidate vCPU

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**

- Boost preempted vCPUs where they have received an IPI [2019]

Insights

Lost opportunity

- Spinlock still incurs **TLB shootdowns**
- **vCPUs in user mode** also need to ack when **TLB shutdown happens**
- **No need** to boost **IPI-receivers** when **lock-holder** **preemption** happens

How to select candidate vCPU

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**
- Boost preempted vCPUs where they have received an IPI [2019]

Insights

Optimization against **TLB shutdown**

Overboost

- Spinlock still in
- **vCPUs in user mode** need to ack when **TLB shutdown happens**
- **No need to boost IPI-receivers when lock-holder preemption happens**

Mitigation: Strict Boost

- PLE-ed vCPU sleeps (NOT boost any vCPUs) [2009]
 - Directed yield [2011]
 - Less prioritize recently PLE-ed vCPUs [2012]
 - Boost only preempted vCPUs [2013]
 - Boost only vCPUs in kernel mode [2017]
- } Optimizations against **spinlock**
- **B** Introduce new strategies for candidate vCPU selection [2019]
 - **Boost IPI-receiver vCPUs in user mode**
 - **Not boost halted vCPUs if the yielded vCPU has NOT sent an IPI to it**

Why PLE events occur continuously?

- ~~Lost opportunity and Overboost~~ **Strict Boost**
- Scheduler mismatch

Problem: Scheduler Mismatch

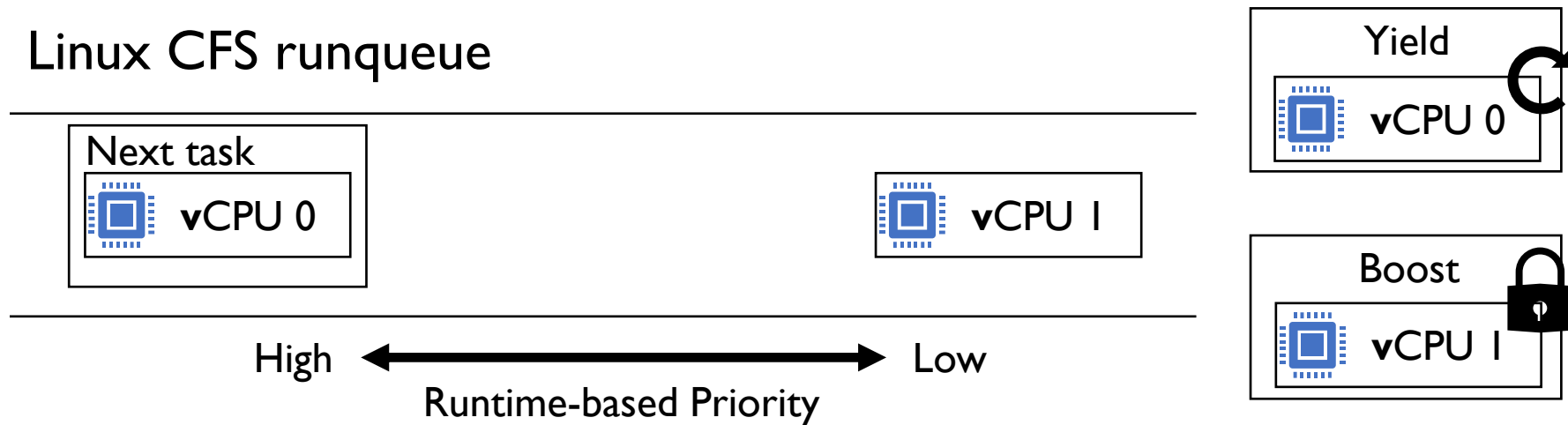
- Linux CFS does **not distinguish** between vCPUs and other threads
 - KVM makes request to Linux CFS for boosting vCPU
 - But Linux CFS always keeps fairness between vCPUs

```
if (cfs_rq->next && wakeup_preempt_entity(cfs_rq->next, left) < 1) {  
    /*  
     * Someone really wants this to run. If it's not unfair, run it.  
     */  
    se = cfs_rq->next;
```

in kernel/sched/fair.c

Case Study: Scheduler Mismatch

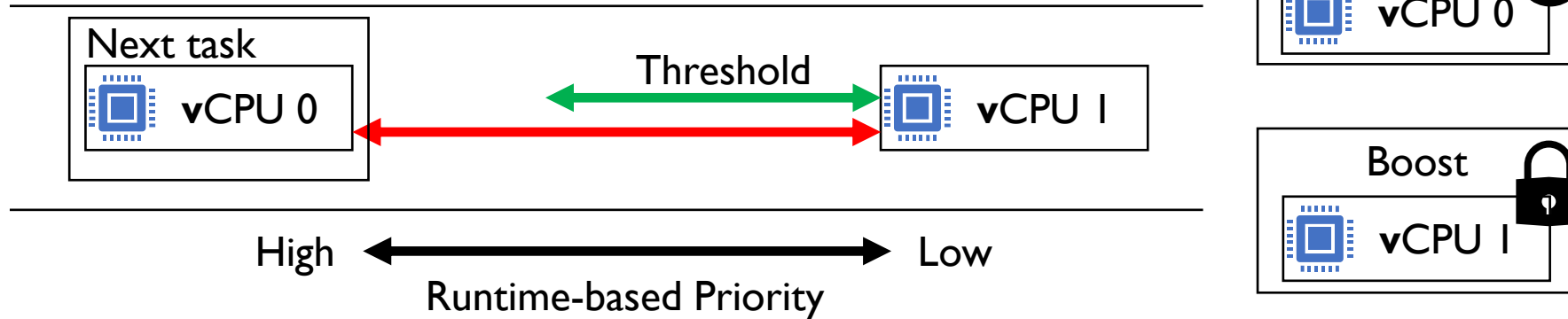
Linux CFS runqueue



1. Picks the highest priority task: vCPU 0
2. $cputime_{vCPU1} - cputime_{vCPU0} > threshold$
3. Decides not to yield vCPU 0 and not to boost vCPU 1
4. vCPU 0 triggers PLE again because vCPU 1 still has the lock

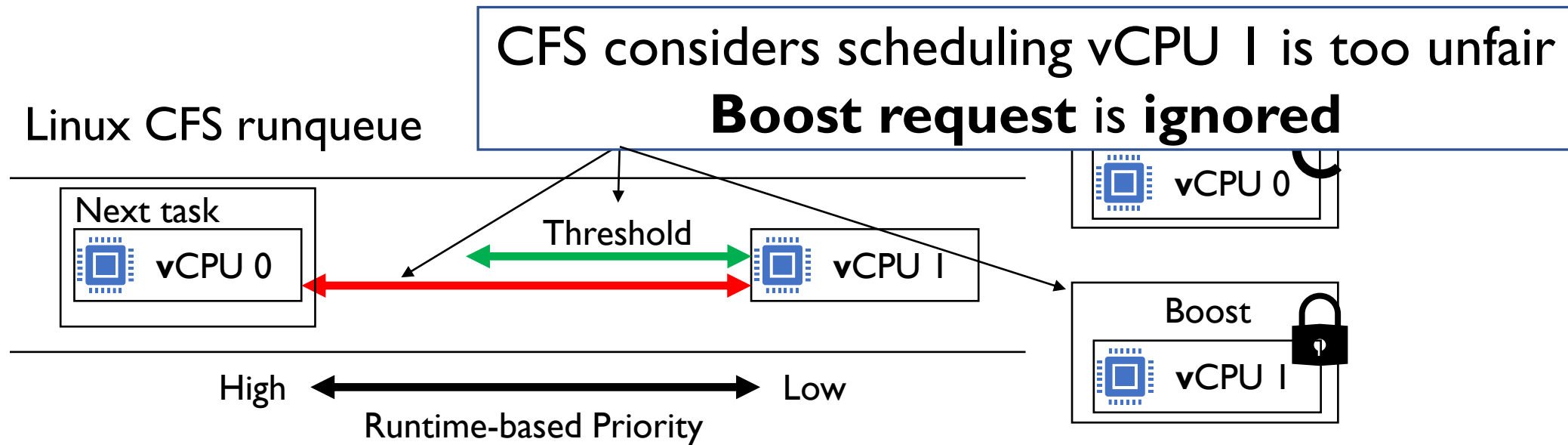
Case Study: Scheduler Mismatch

Linux CFS runqueue



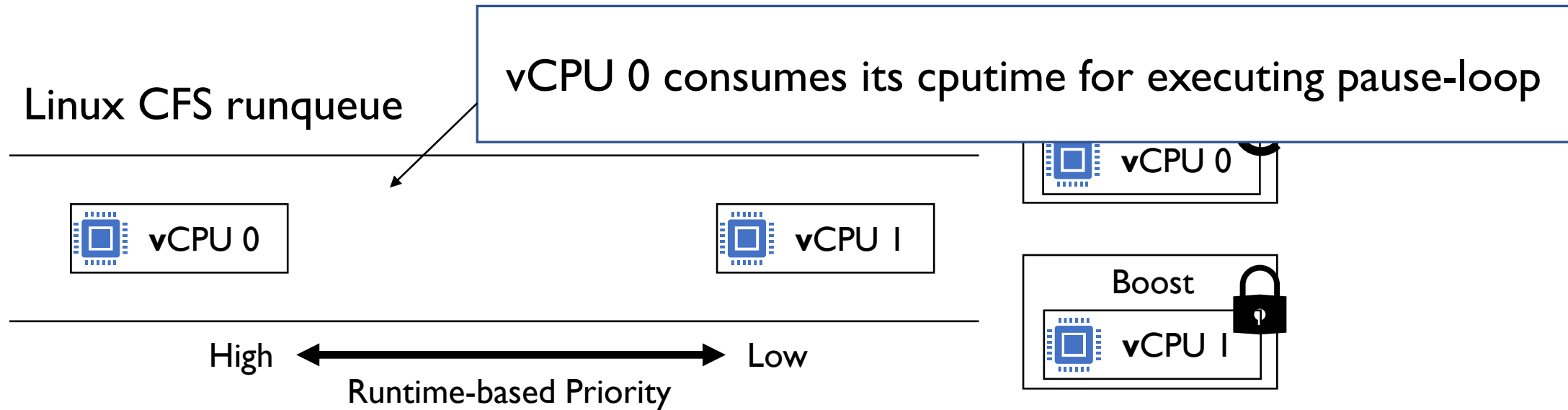
1. Picks the highest priority task: vCPU 0
2. $cputime_{vCPU1} - cputime_{vCPU0} > threshold$
3. Decides not to yield vCPU 0 and not to boost vCPU 1
4. vCPU 0 triggers PLE again because vCPU 1 still has the lock

Case Study: Scheduler Mismatch



1. Picks the highest priority task: vCPU 0
2. $cputime_{vCPU1} - cputime_{vCPU0} > threshold$
3. Decides not to yield vCPU 0 and not to boost vCPU 1
4. vCPU 0 triggers PLE again because vCPU 1 still has the lock

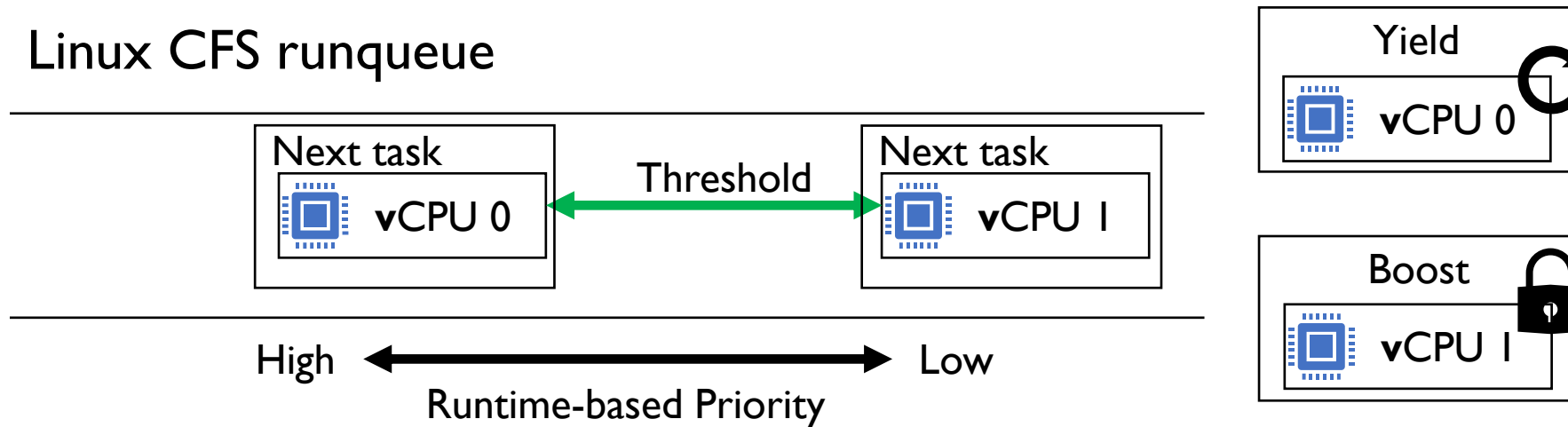
Case Study: Scheduler Mismatch



1. Picks the highest priority task: vCPU 0
2. $cputime_{vCPU1} - cputime_{vCPU0} > threshold$
3. Decides not to yield vCPU 0 and not to boost vCPU 1
4. vCPU 0 triggers PLE again because vCPU 1 still has the lock

Case Study: Scheduler Mismatch

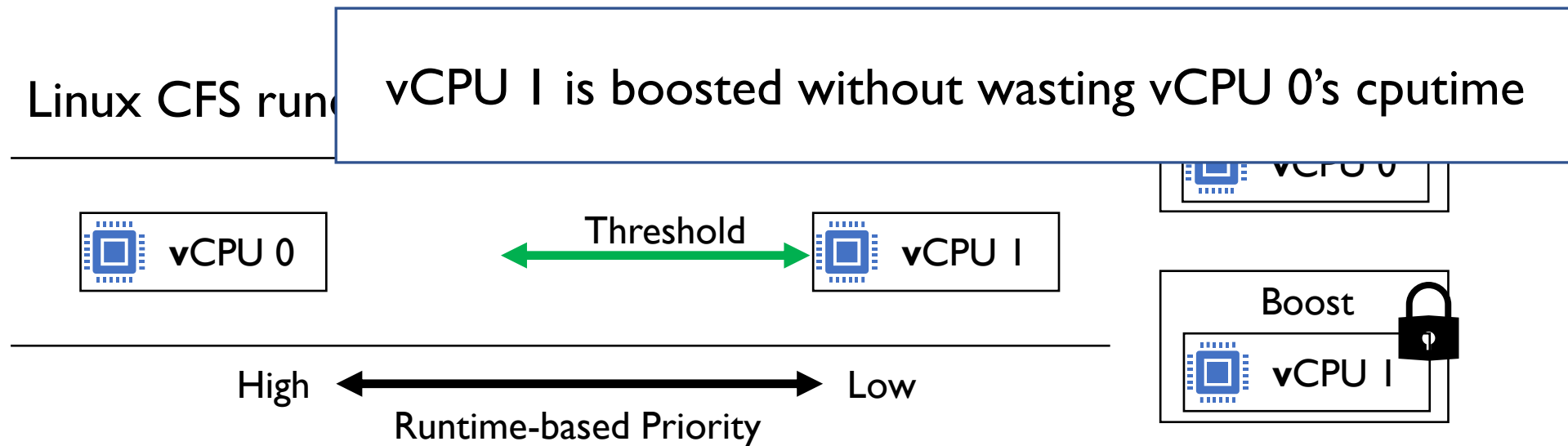
Linux CFS runqueue



1. Picks the highest priority task: vCPU 0
2. vCPU 1 is boosted eventually because $cputime_{vCPU1} - cputime_{vCPU0} < threshold$

Mitigation: Debooster

- Debooster makes CFS not to hesitate to boost another vCPU instead of vCPU which exits due to PLE
 - by lowering PLE-ed vCPU priority



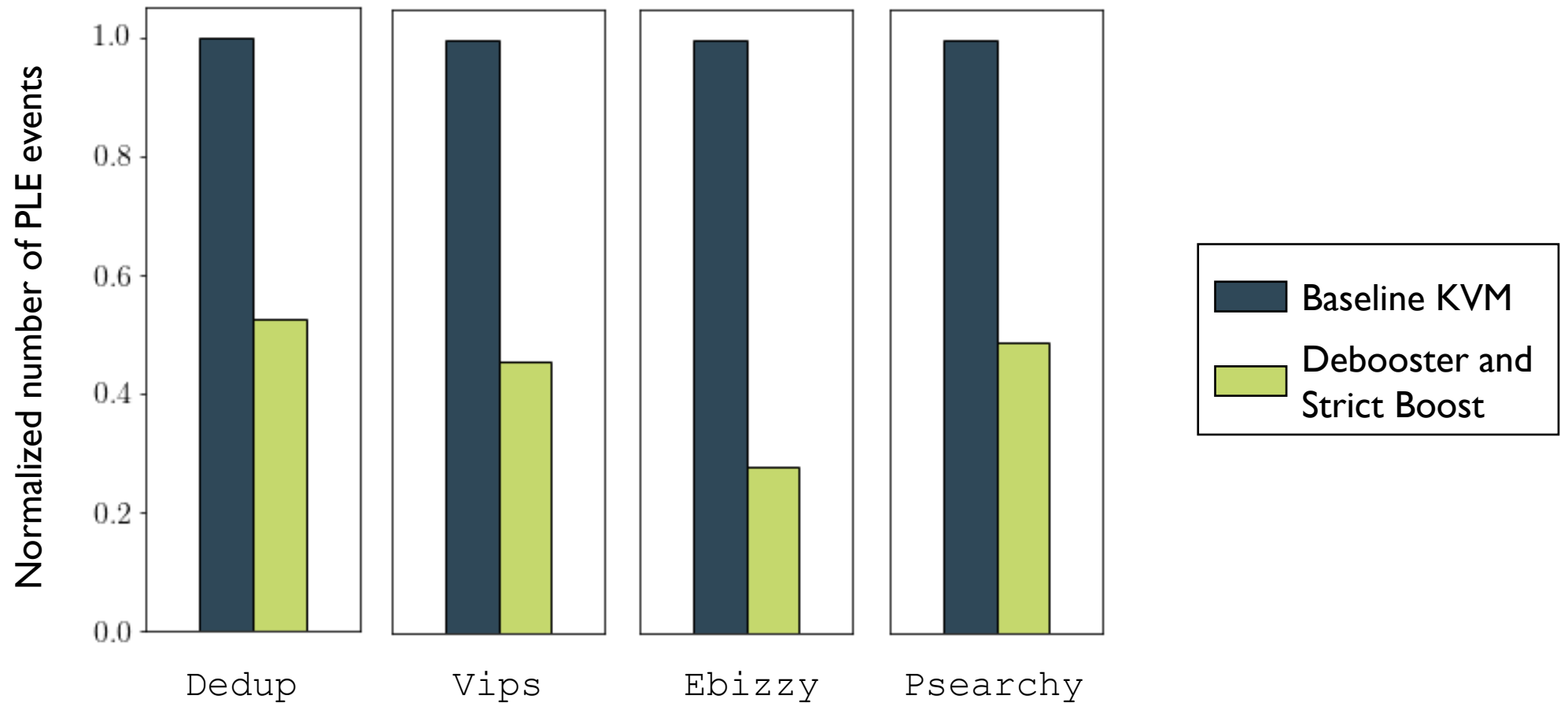
Implementation

- Modified Linux KVM 5.6.0
- Our mitigations require 4l LoC modification
 - Debooster: `yield_to_task()` interface in CFS
 - Strict Boost: vCPU candidate selection and IPI handler in KVM

Experimental Setup

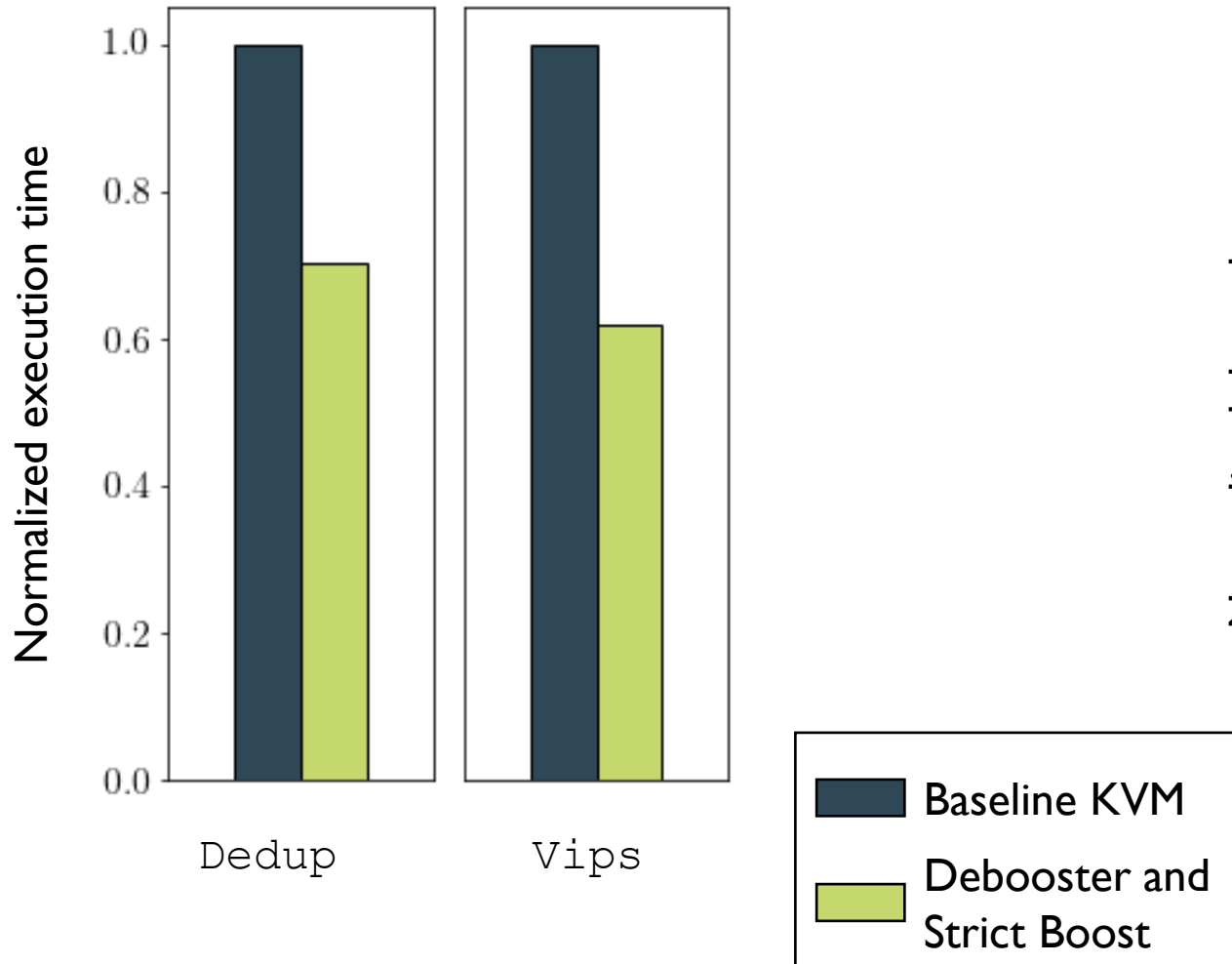
- Testbed
 - 8-core (Intel Xeon)
 - 2 VMs with 8 vCPU for each
- Benchmarks
 - `dedup` and `vips` from `parsec`
 - `psearchy` from `mosbench`
 - `ebizzy`
 - `swaptions`: CPU-intensive workload co-runner in another VM

Evaluation: Reduction of PLE Occurrences

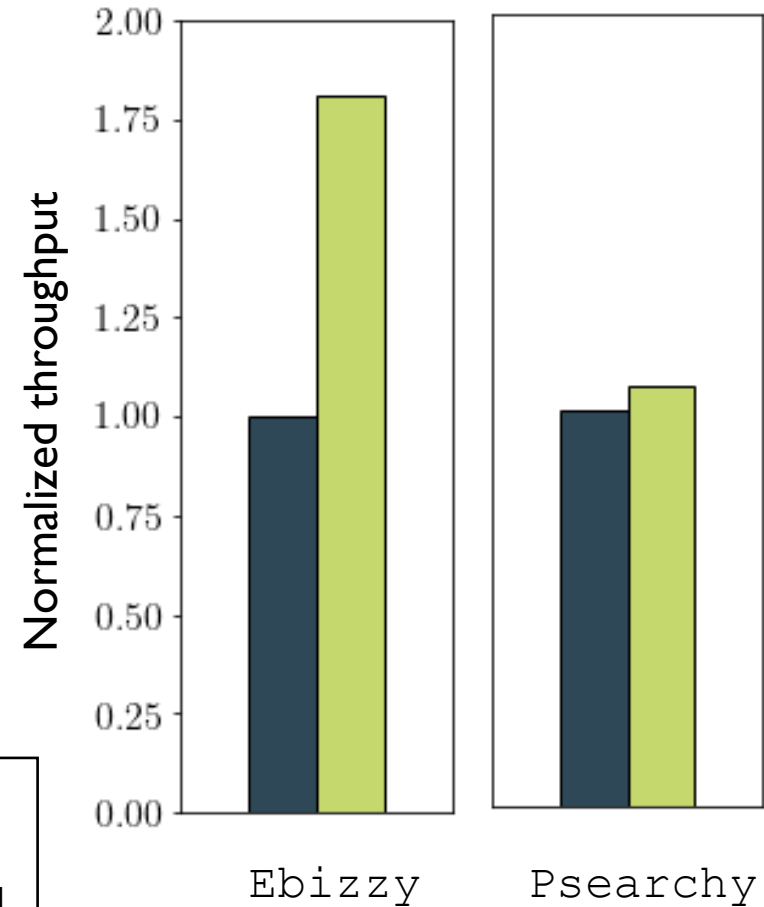


Evaluation: Performance Improvement

Execution time is **reduced by 40 %**

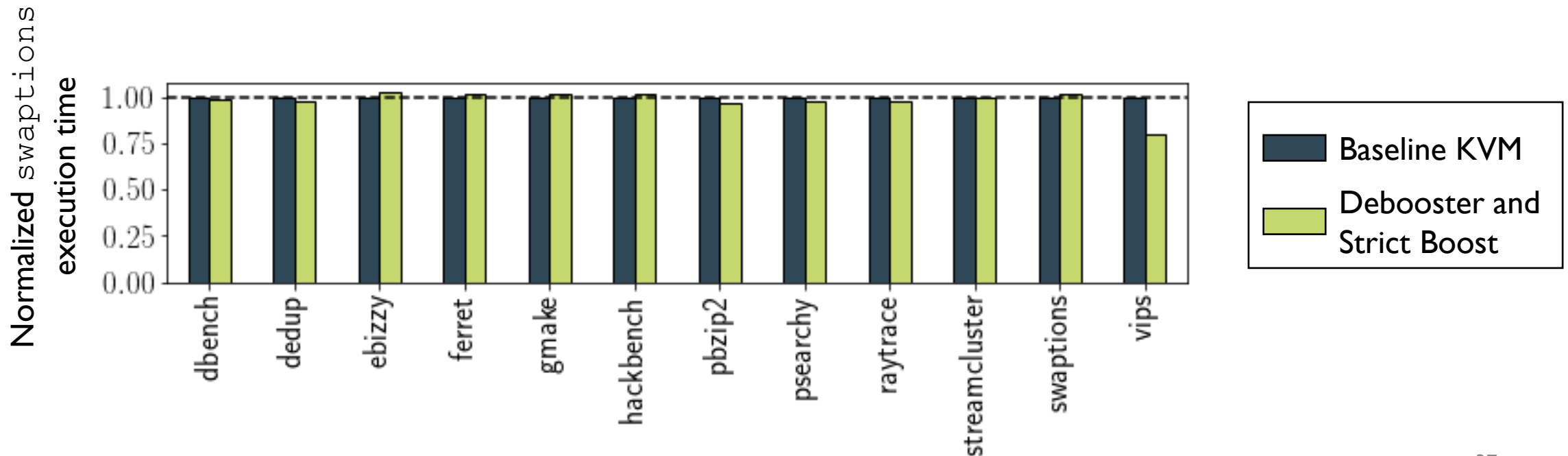


Throughput is **improved by 75 %**



Evaluation: System Fairness

- Co-runner's performance degradation cannot be seen
- Mitigations do not raise the priority of the boosted vCPU



Conclusion

- Oversubscribing vCPUs incurs excessive spinning
- Pause-Loop-Exiting (PLE) unfortunately does not fix it
- due to the semantic gap between
 - KVM and guest VMs
 - KVM and Linux scheduler
- Introduced mitigations against identified problems improve apps throughput by up to 75 %
 - **Strict Boost** against **lost opportunity** and **overboost**
 - **Debooster** against **scheduler mismatch**
- Problem investigated by the KVM community
 - <https://lore.kernel.org/kvm/20210421150831.60133-1-kentaishiguro@sslab.ics.keio.ac.jp/>
 - <https://lore.kernel.org/kvm/1618542490-14756-1-git-send-email-wanpengli@tencent.com/>
- Ishiguro, K., Yasuno, N., Aublin, P. L., & Kono, K. (2021, April). Mitigating excessive vCPU spinning in VM-agnostic KVM (VEE '21)

Thanks & QA

Mitigating Excessive Pause-Loop-Exiting in VM-Agnostic KVM

