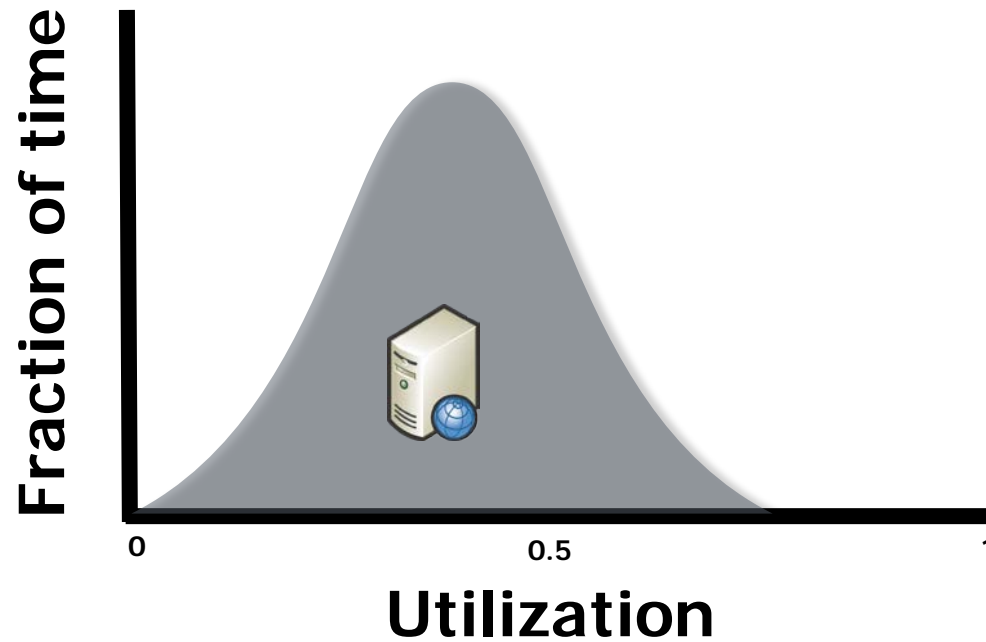


Accelerating Critical OS Services in Virtualized Systems with Flexible Micro-sliced Cores

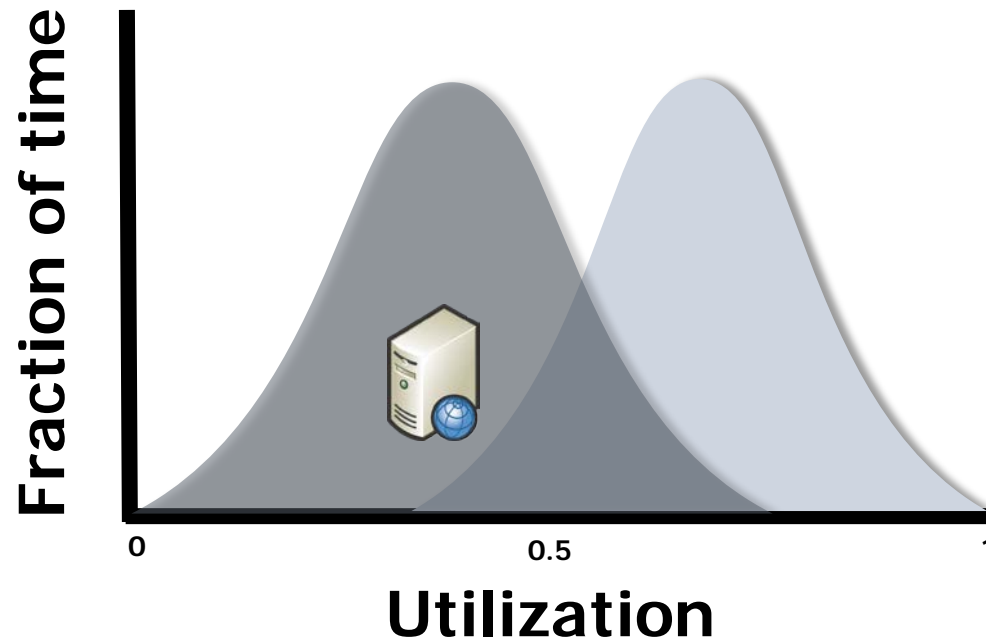
Jeongseob Ahn^{*}, Chang Hyun Park[‡], Taekyung Heo[‡], Jaehyuk Huh[‡]



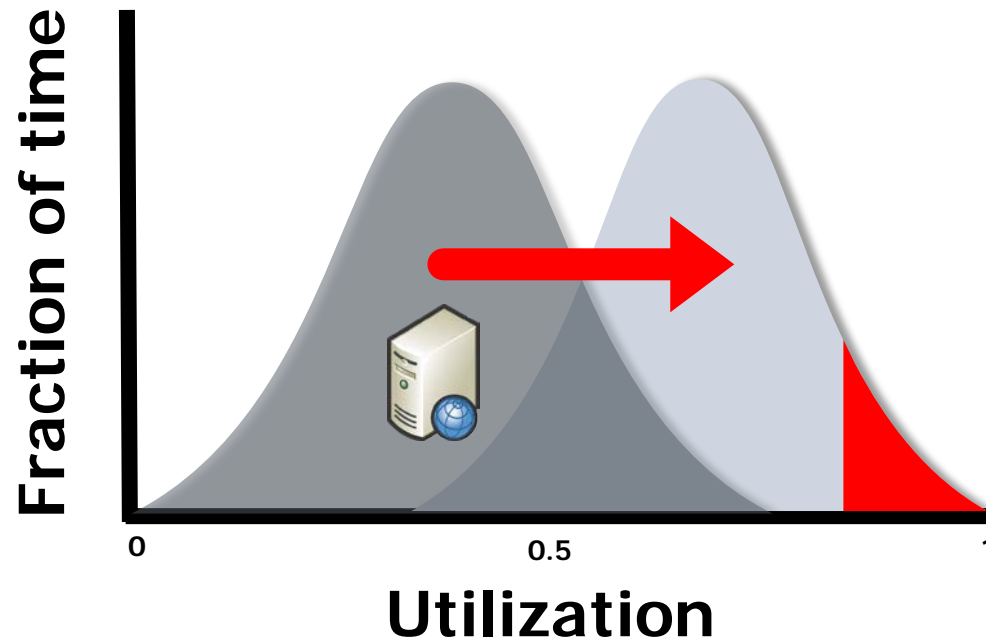
Challenge of Server Consolidation



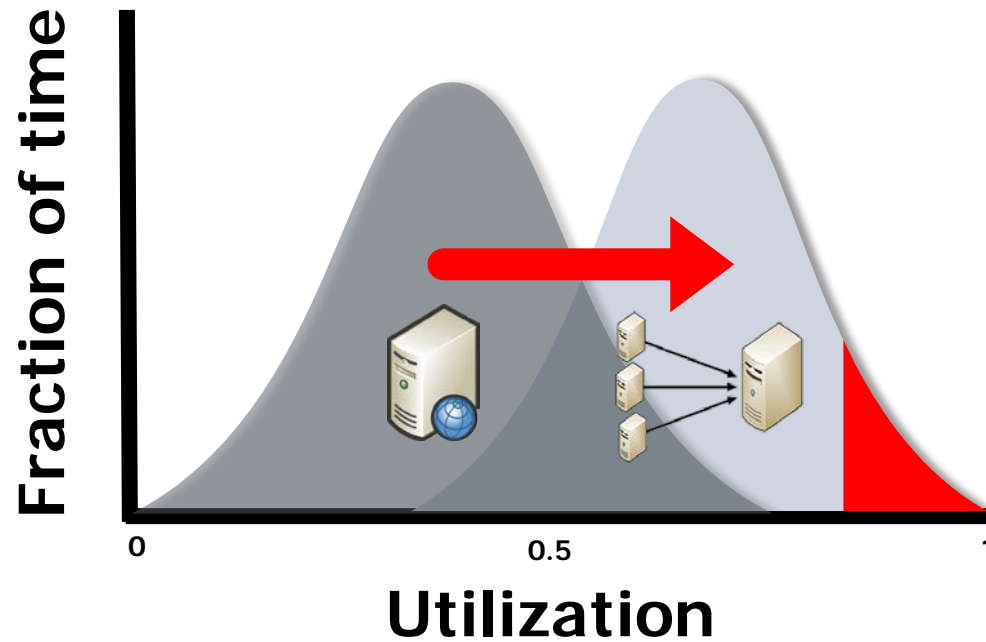
Challenge of Server Consolidation



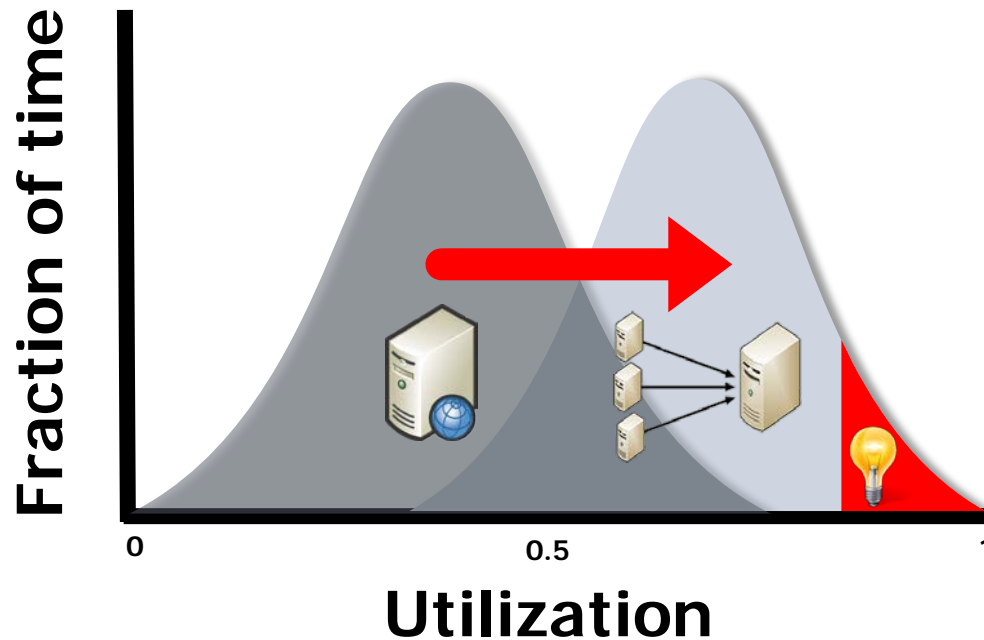
Challenge of Server Consolidation



Challenge of Server Consolidation

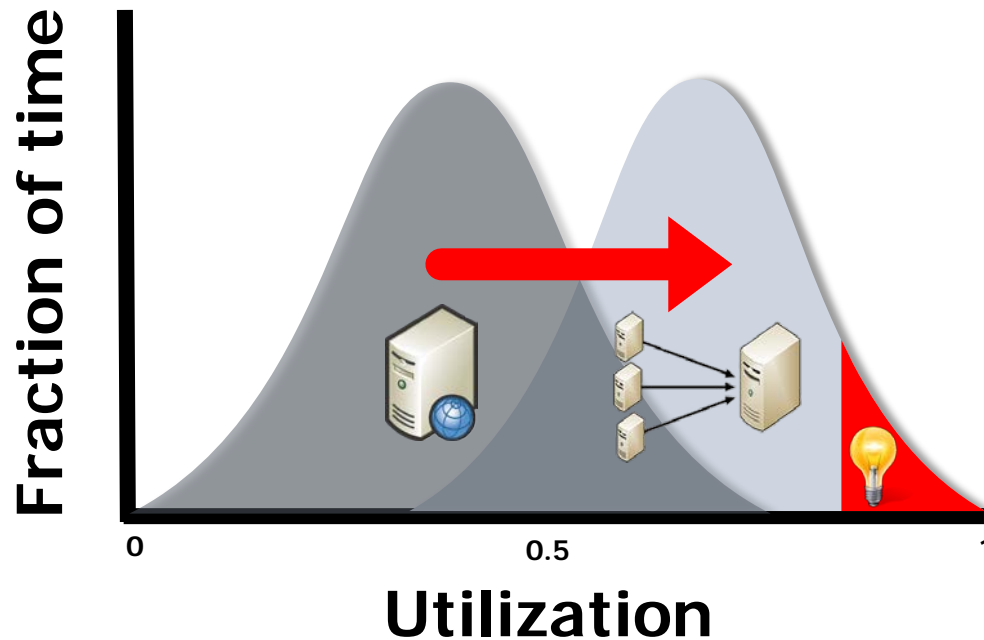


Challenge of Server Consolidation



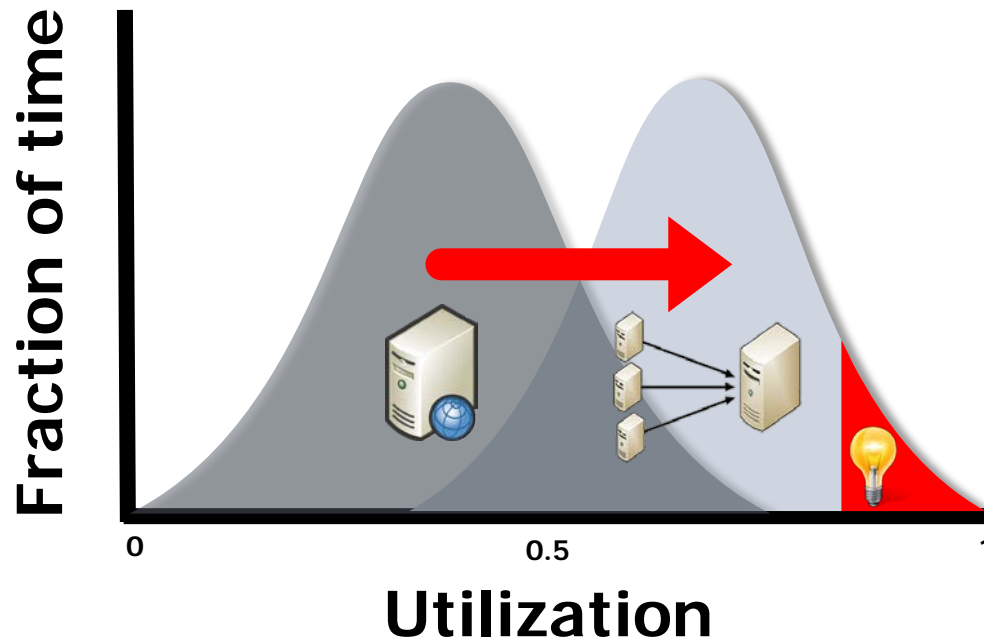
Consolidation improves system utilization

Challenge of Server Consolidation



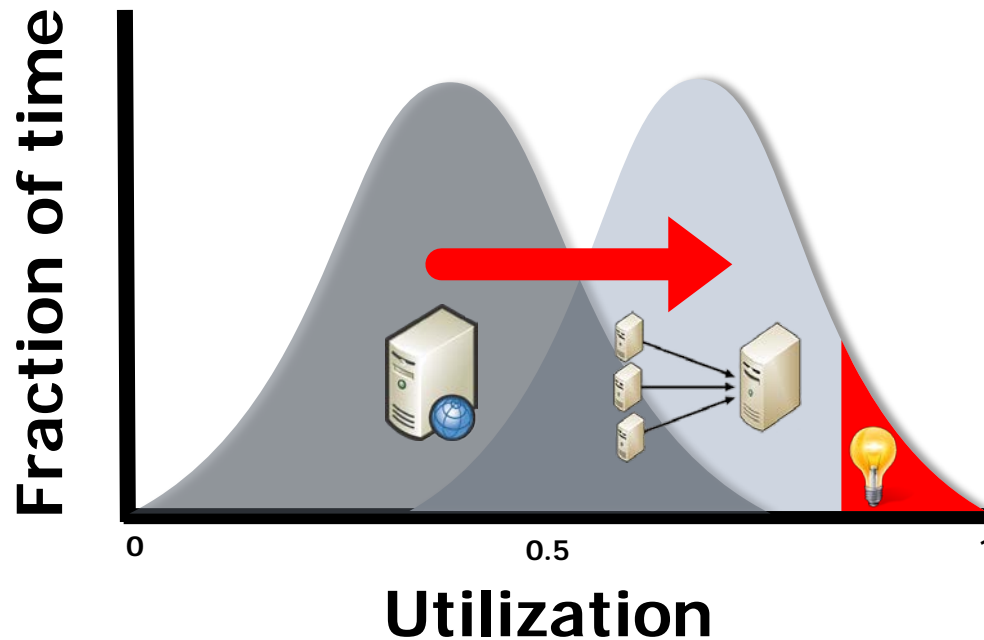
Consolidation improves system utilization

Challenge of Server Consolidation



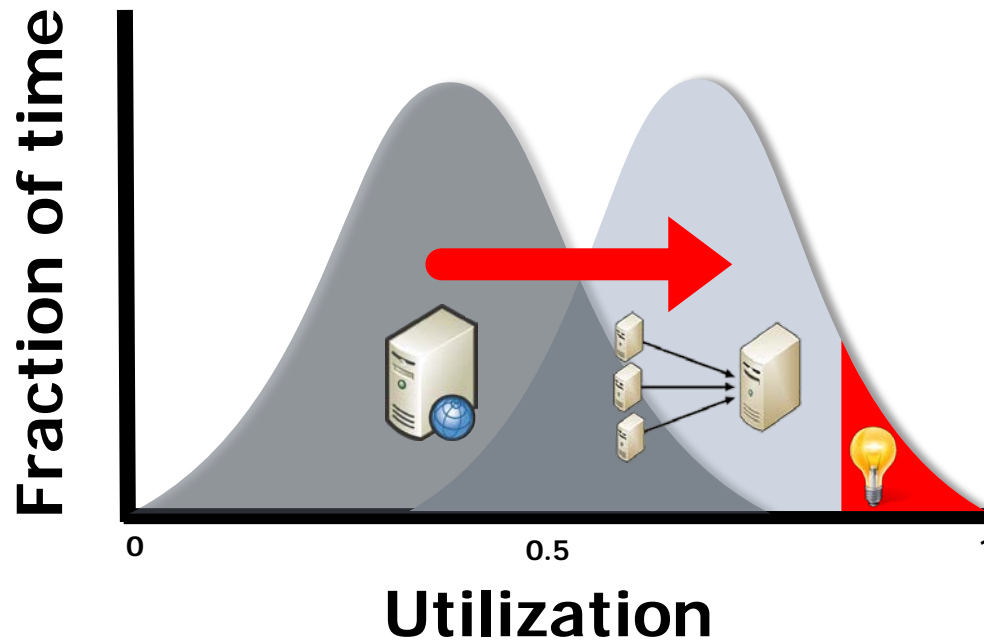
Consolidation improves system utilization

Challenge of Server Consolidation



Consolidation improves system utilization
However, resources are contended

Challenge of Server Consolidation

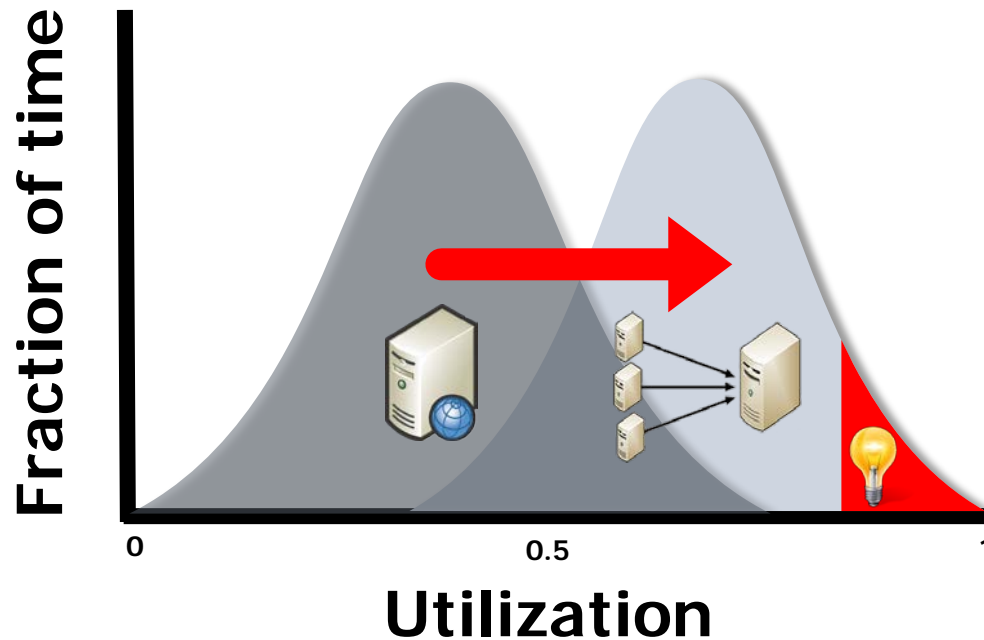


Consolidation improves system utilization



However, resources are contended

Challenge of Server Consolidation

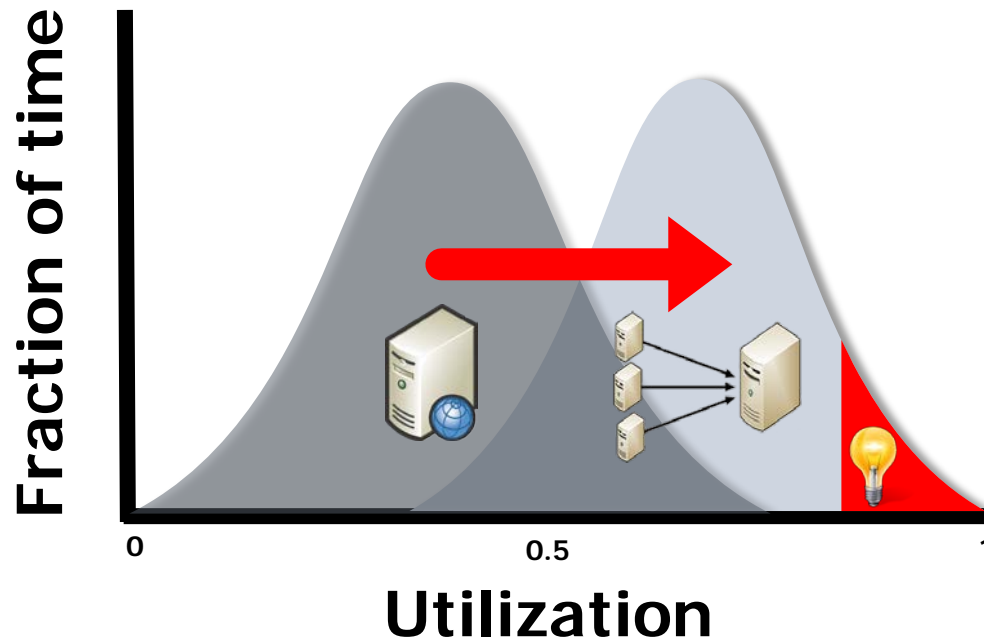


Consolidation improves system utilization



However, resources are contended

Challenge of Server Consolidation



Consolidation improves system utilization



However, resources are contended

Challenge of Server Consolidation



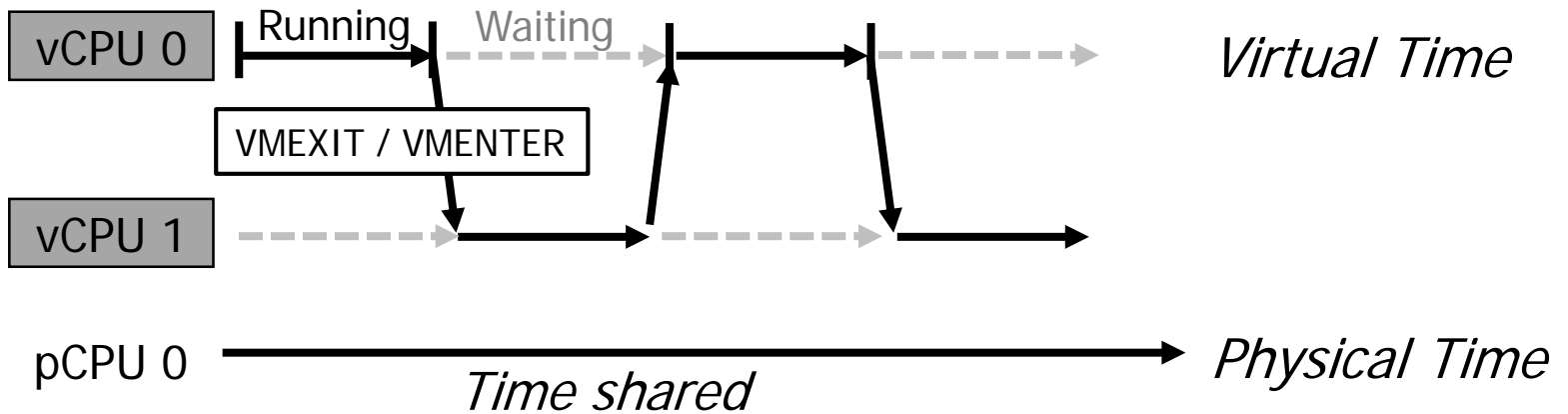
Consolidation improves system utilization



However, resources are contended

So, What Can Happen?

- Virtual time discontinuity



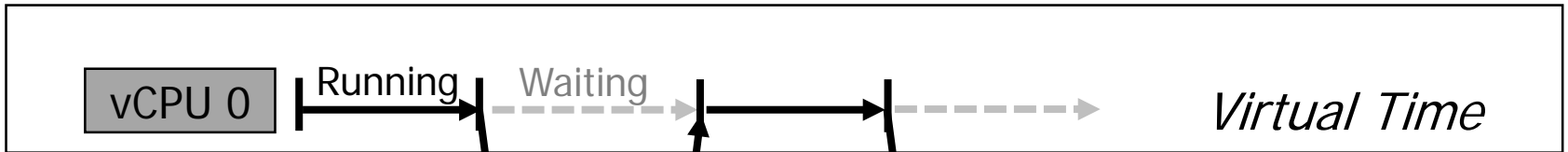
So, What Can Happen?

- Virtual time discontinuity



So, What Can Happen?

- Virtual time discontinuity



① Spinlock waiting time (gmake)

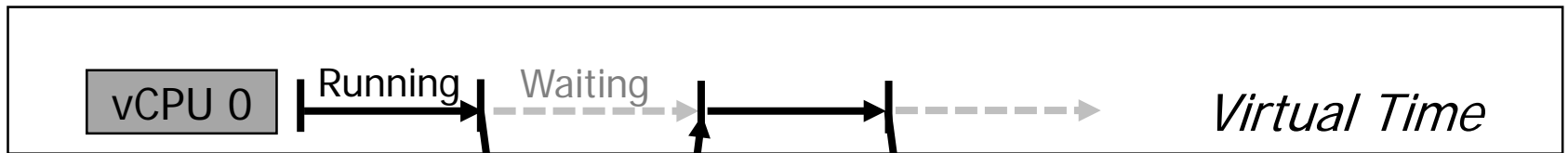
Kernel Component	Avg. waiting time (μsec)	
	solo	co-run*
Page reclaim	1.03	420.13
Page allocator	3.42	1,053.26
Dentry	2.93	1,298.87
Runqueue	1.22	256.07

* Concurrently running with Swaptions of PARSEC

Processing time is amplified

So, What Can Happen?

- Virtual time discontinuity



① Spinlock waiting time (gmake)

Kernel Component	Avg. waiting time (μsec)	
	solo	co-run*
Page reclaim	1.03	420.13
Page allocator	3.42	1,053.26
Dentry	2.93	1,298.87
Runqueue	1.22	256.07

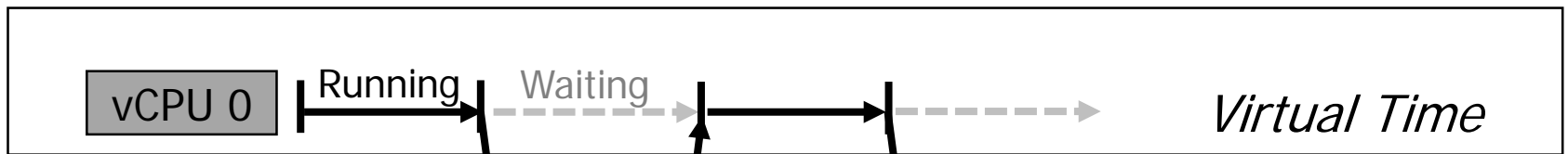
		Avg.	Min.	Max.
dedup	solo	28	5	1927
	co-run*	6,354	7	74915
vips	solo	55	5	2052
	co-run*	14,928	17	121548

* Concurrently running with Swaptions of PARSEC

Processing time is amplified

So, What Can Happen?

- Virtual time discontinuity



① Spinlock waiting time (gmake)

Kernel Component	Avg. waiting time (μsec)	
	solo	co-run*
Page reclaim	1.03	420.13
Page allocator	3.42	1,053.26
Dentry	2.93	1,298.87
Runqueue	1.22	256.07

* Concurrently running with Swaptions of PARSEC

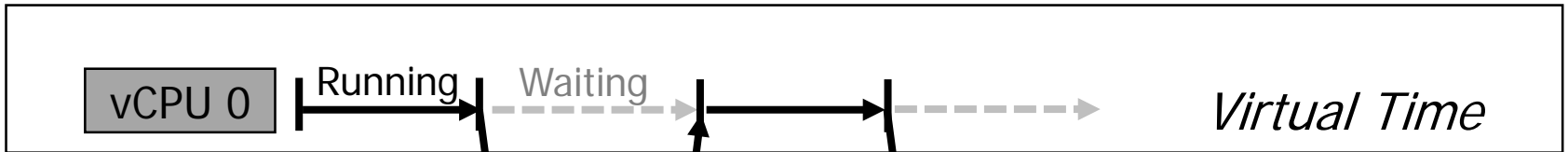
② TLB synchronization latency (μsec)

		Avg.	Min.	Max.
dedup	solo	28	5	1927
	co-run*	6,354	7	74915
vips	solo	55	5	2052
	co-run*	14,928	17	121548

Processing time is amplified

So, What Can Happen?

- Virtual time discontinuity



① Spinlock waiting time (gmake)

Kernel Component	Avg. waiting time (μsec)	
	solo	co-run*
Page reclaim	1.03	420.13
Page allocator	3.42	1,053.26
Dentry	2.93	1,298.87
Runqueue	1.22	256.07

* Concurrently running with Swaptions of PARSEC

Processing time is amplified

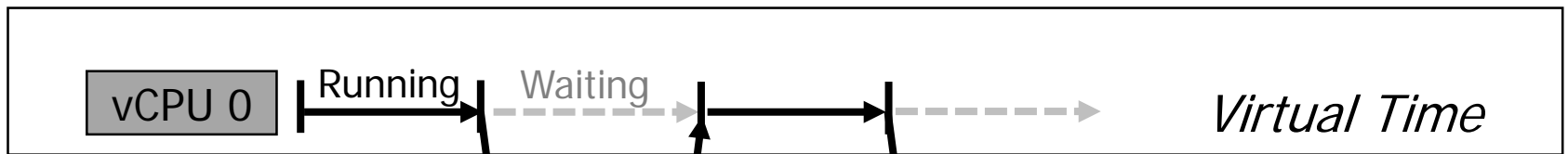
② TLB synchronization latency (μsec)

		Avg.	Min.	Max.
dedup	solo	28	5	1927
	co-run*	6,354	7	74915
vips	solo	55	5	2052
	co-run*	14,928	17	121548

	Jitters (ms)	Throughput (Mbits/sec)
solo	0.0043	936.3
mixed co-run*	9.2507	435.6

So, What Can Happen?

- Virtual time discontinuity



① Spinlock waiting time (gmake)

Kernel Component	Avg. waiting time (μsec)	
	solo	co-run*
Page reclaim	1.03	420.13
Page allocator	3.42	1,053.26
Dentry	2.93	1,298.87
Runqueue	1.22	256.07

* Concurrently running with Swaptions of PARSEC

Processing time is amplified

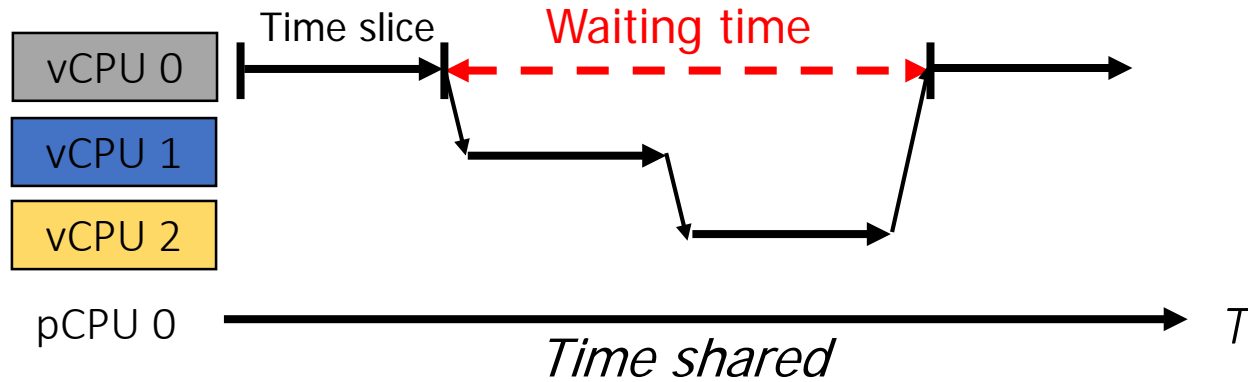
② TLB synchronization latency (μsec)

		Avg.	Min.	Max.
dedup	solo	28	5	1927
	co-run*	6,354	7	74915
vips	solo	55	5	2052
	co-run*	14,928	17	121548

③ I/O latency & throughput (iPerf)

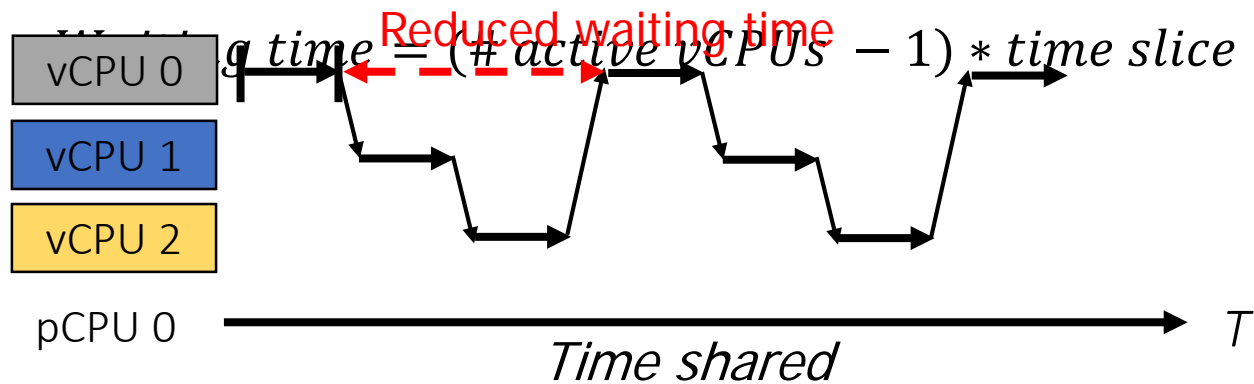
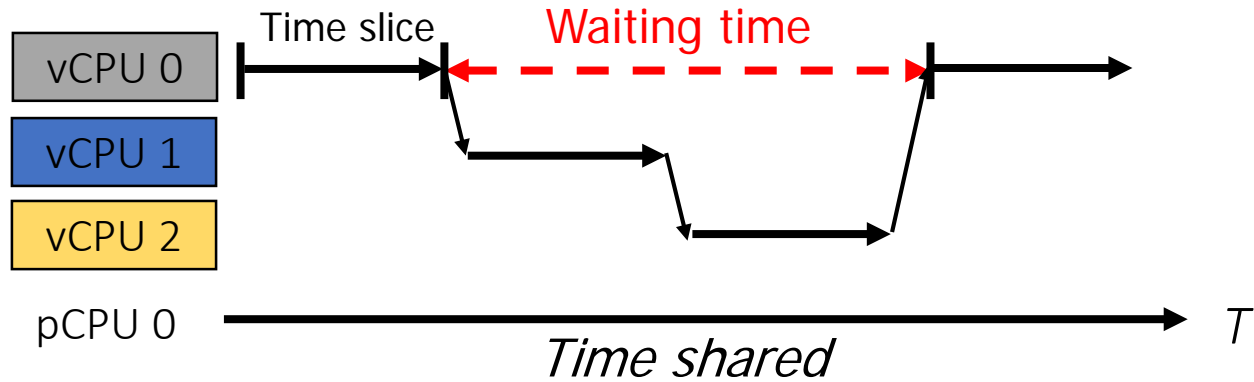
	Jitters (ms)	Throughput (Mbits/sec)
solo	0.0043	936.3
mixed co-run*	9.2507	435.6

How about Shortening Time Slice?

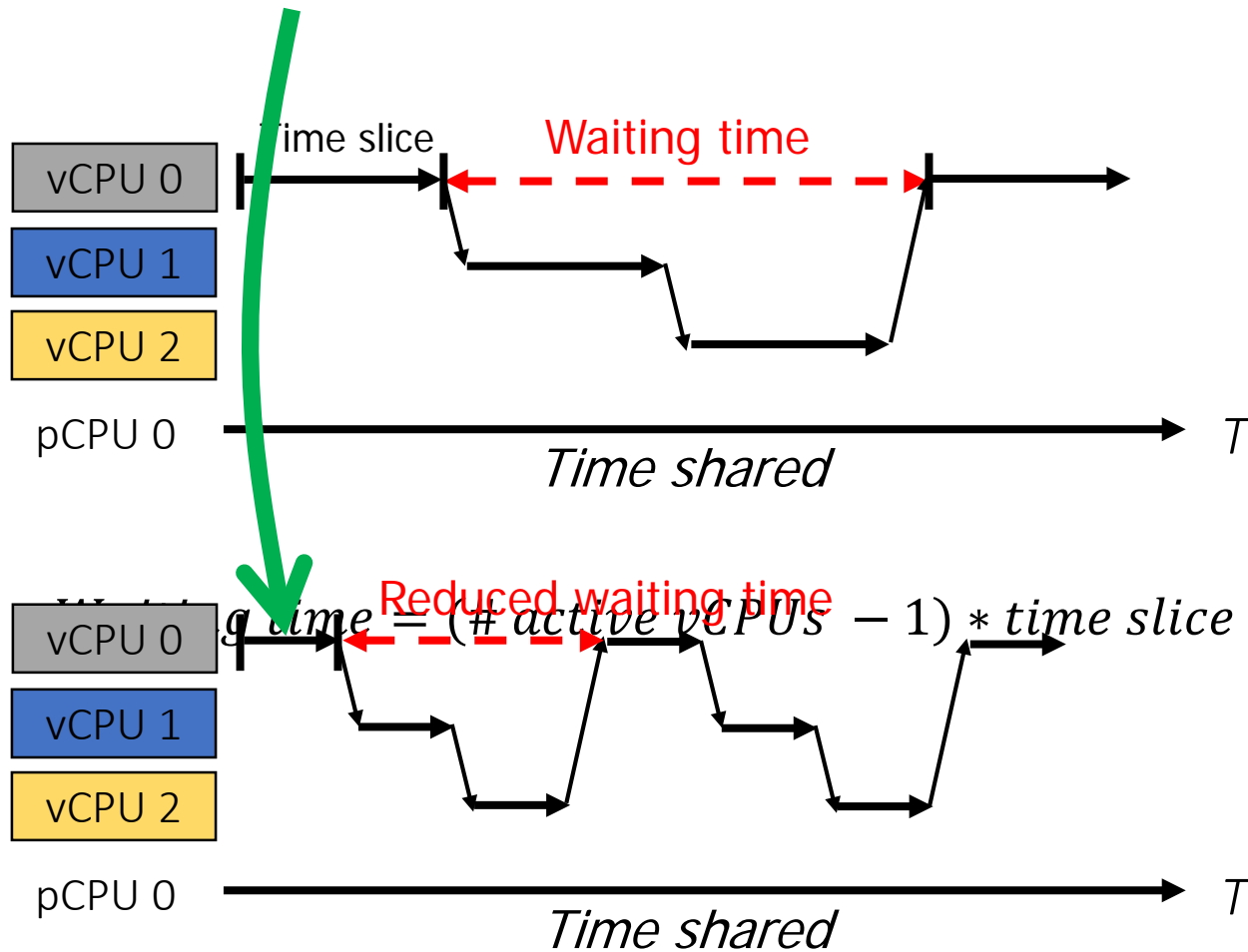


$$\text{Waiting time} = (\# \text{ active vCPUs} - 1) * \text{time slice}$$

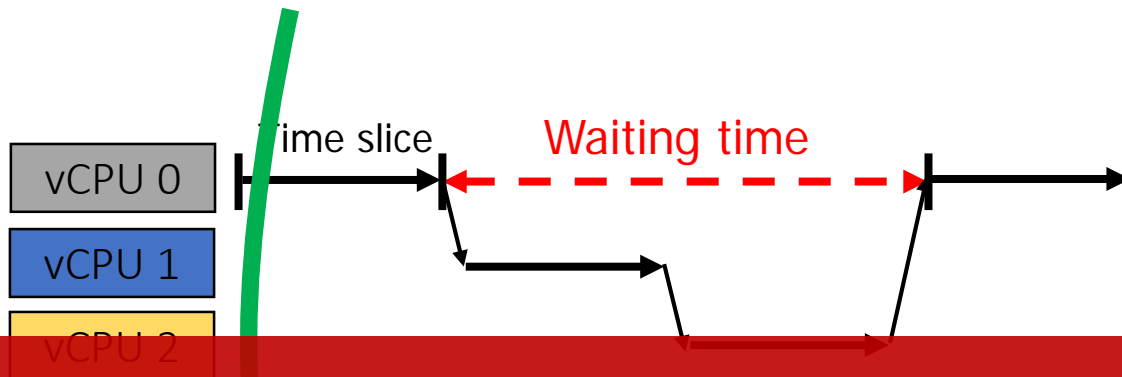
How about Shortening Time Slice?



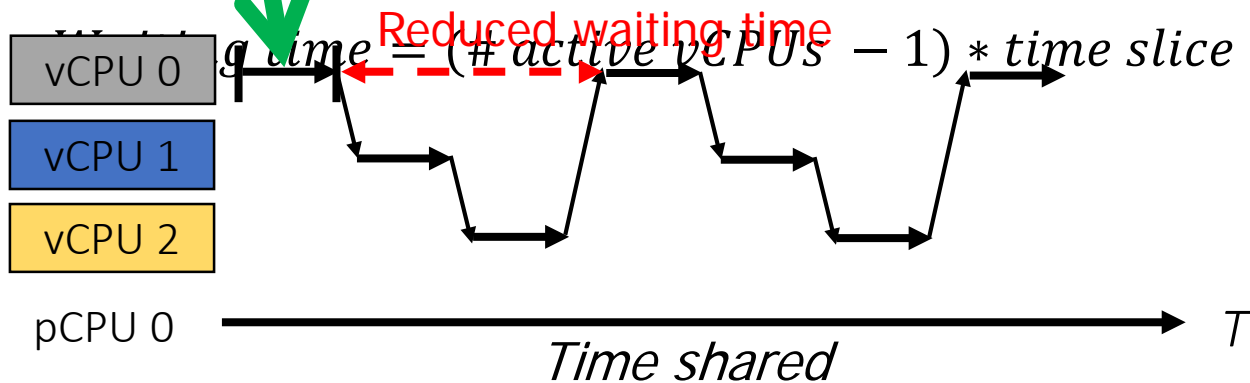
How about Shortening Time Slice?



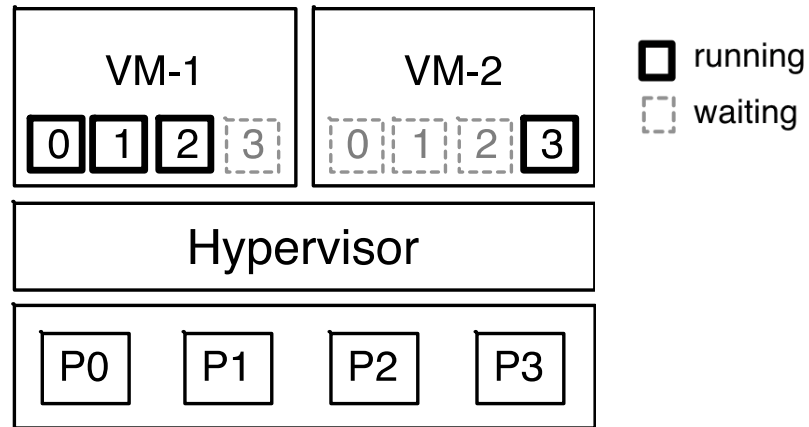
How about Shortening Time Slice?



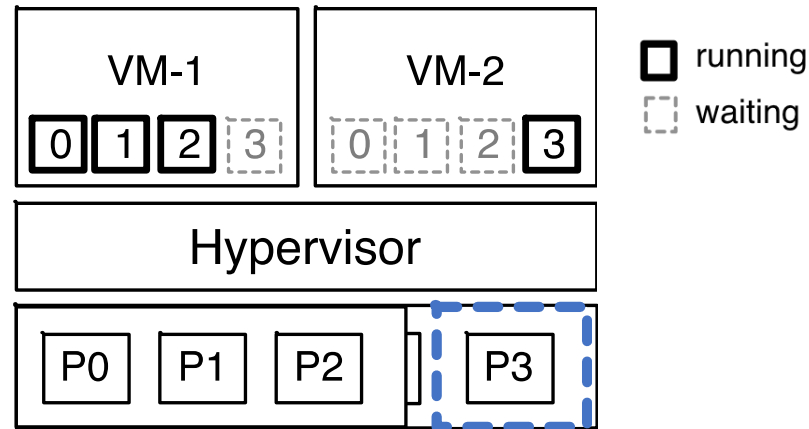
Shortening time slice is very simple and powerful, but the overhead of frequent context switches is significant



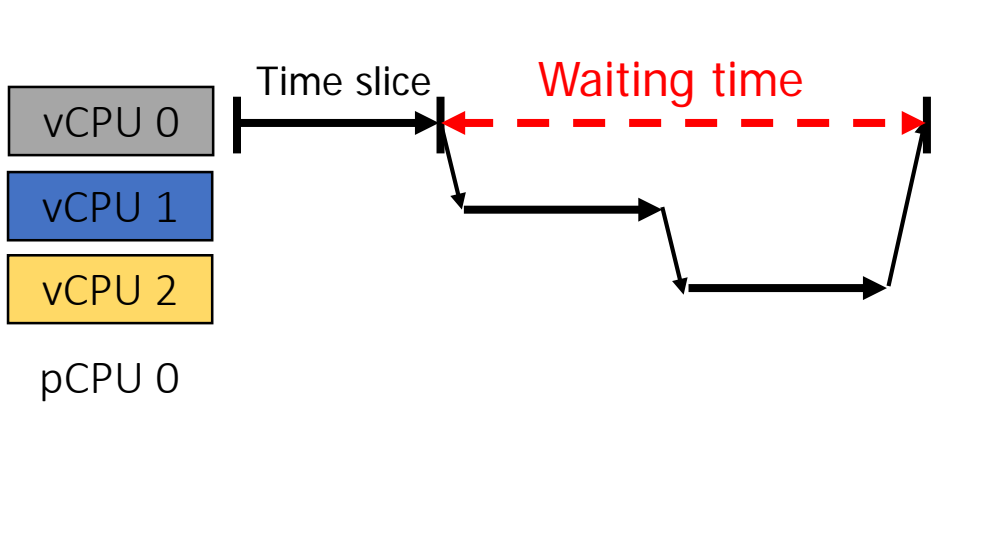
Approach: Dividing CPUs into Two Pools



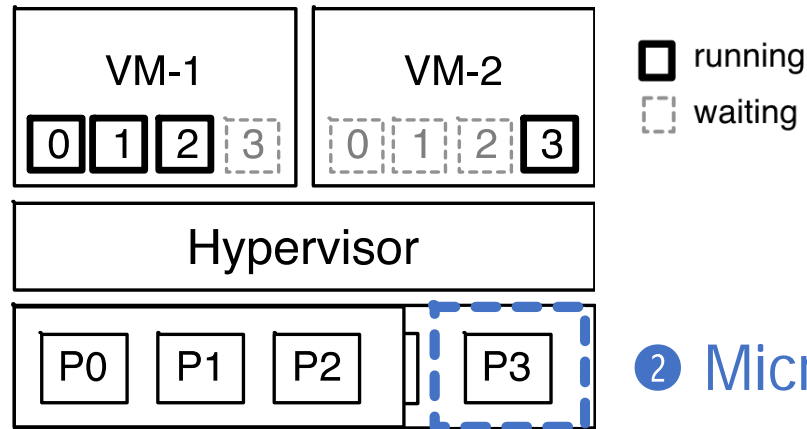
Approach: Dividing CPUs into Two Pools



① Normal pool

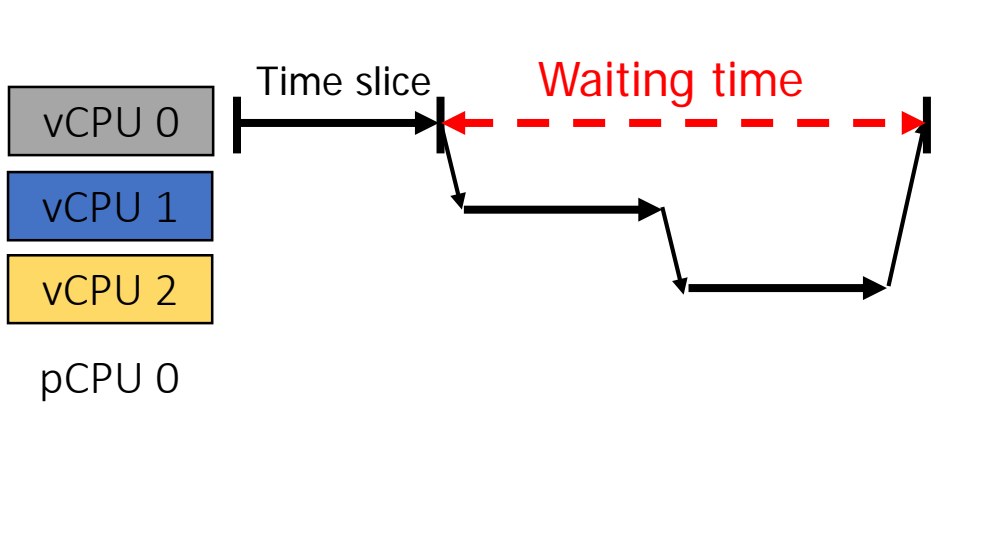


Approach: Dividing CPUs into Two Pools

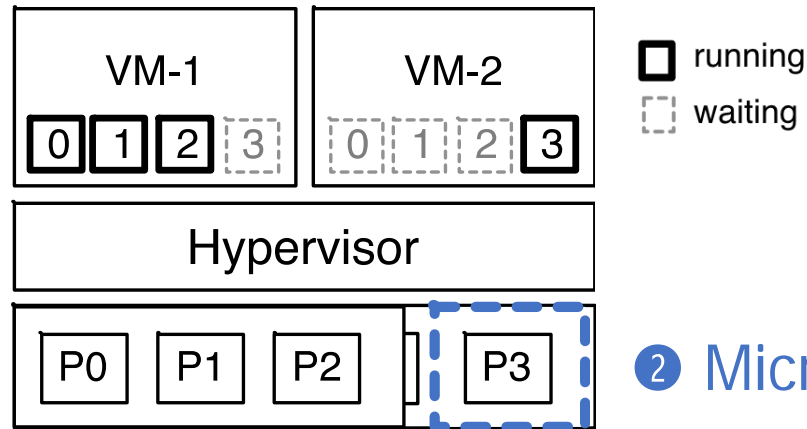


① Normal pool

② Micro-sliced pool

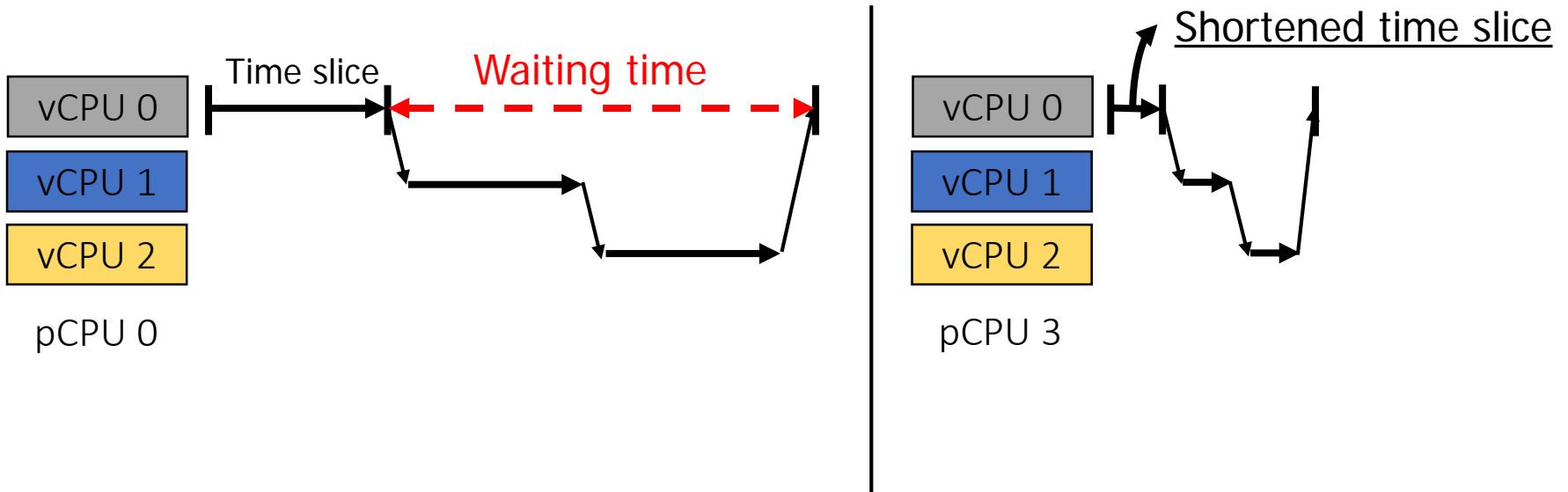


Approach: Dividing CPUs into Two Pools

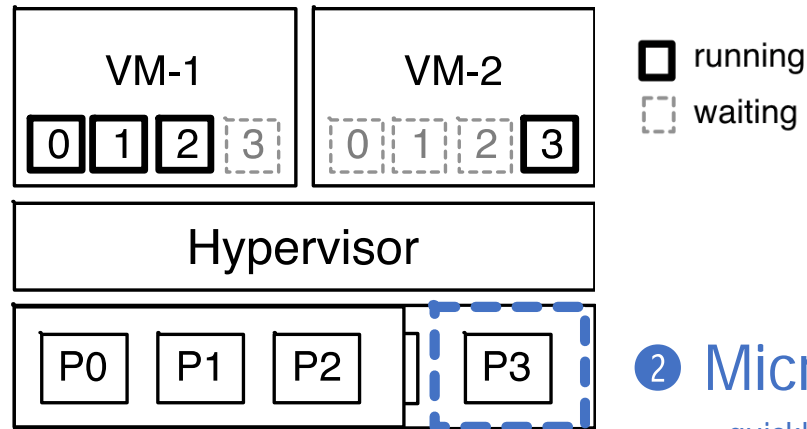


① Normal pool

② Micro-sliced pool

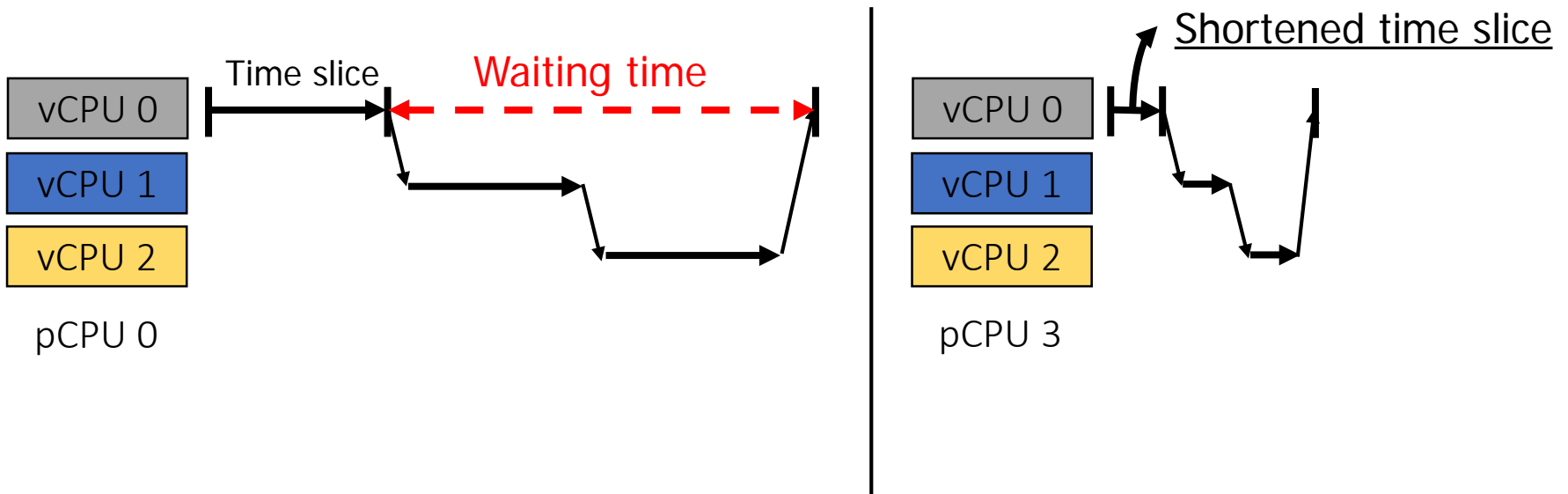


Approach: Dividing CPUs into Two Pools

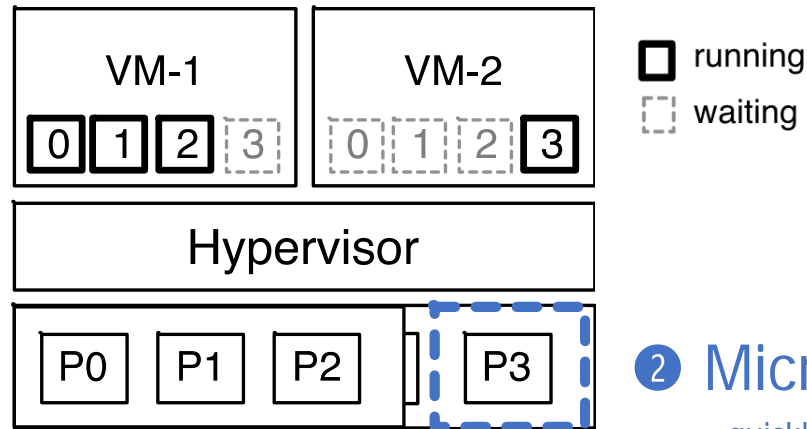


① Normal pool

② Micro-sliced pool
- quickly but briefly schedule vCPUs

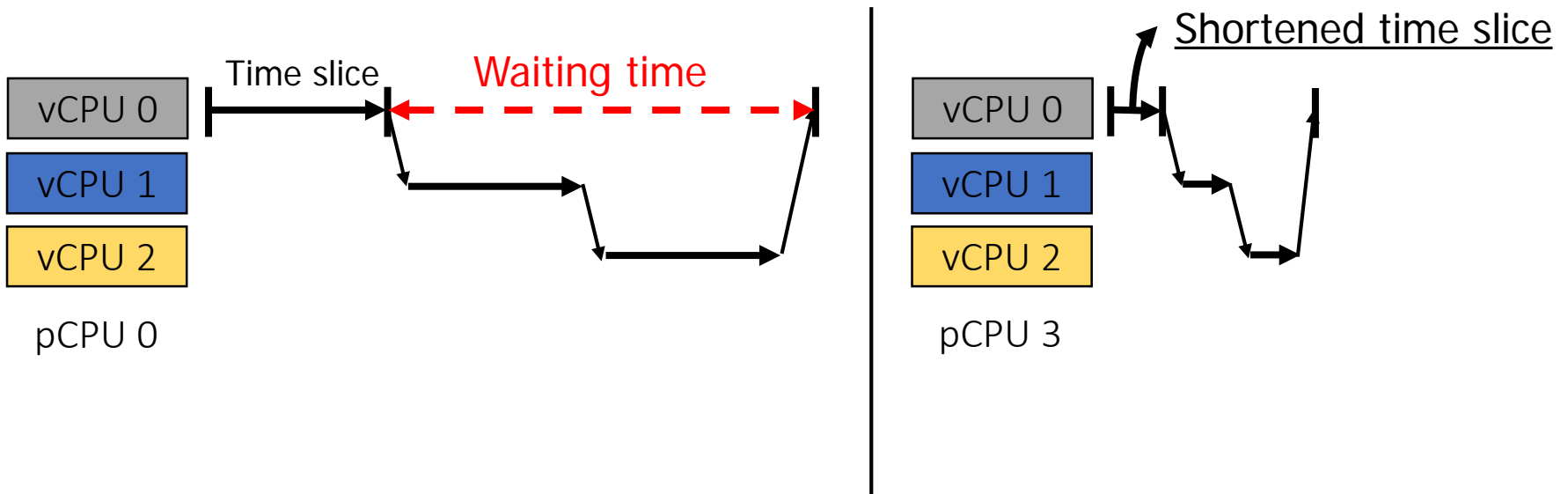


Approach: Dividing CPUs into Two Pools

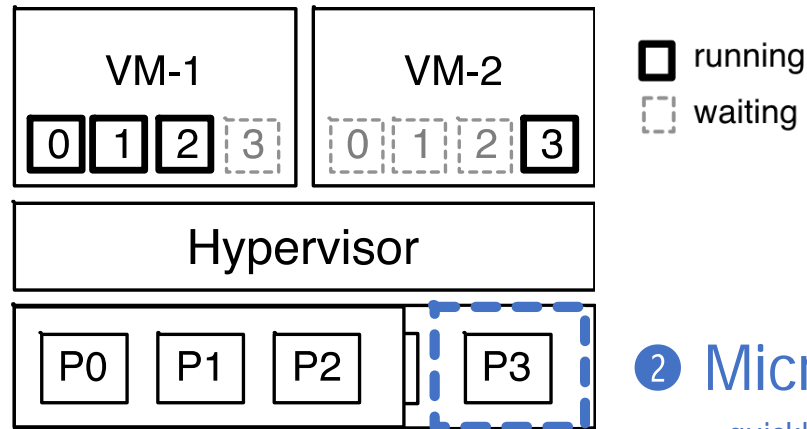


① Normal pool

② Micro-sliced pool
- quickly but briefly schedule vCPUs

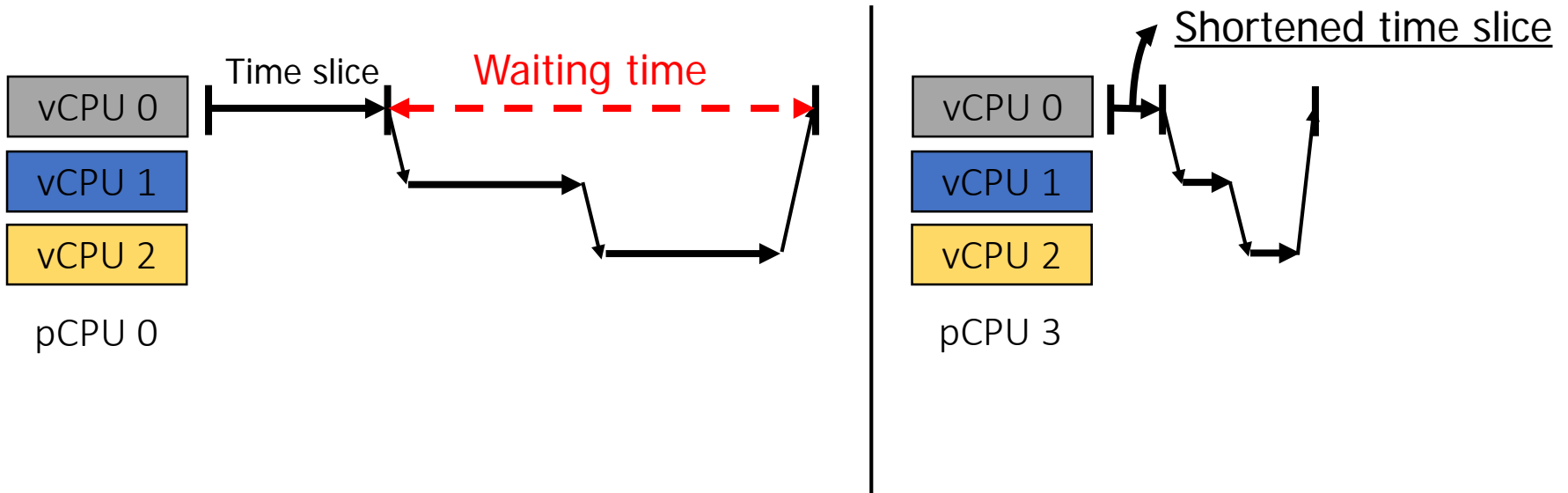


Approach: Dividing CPUs into Two Pools

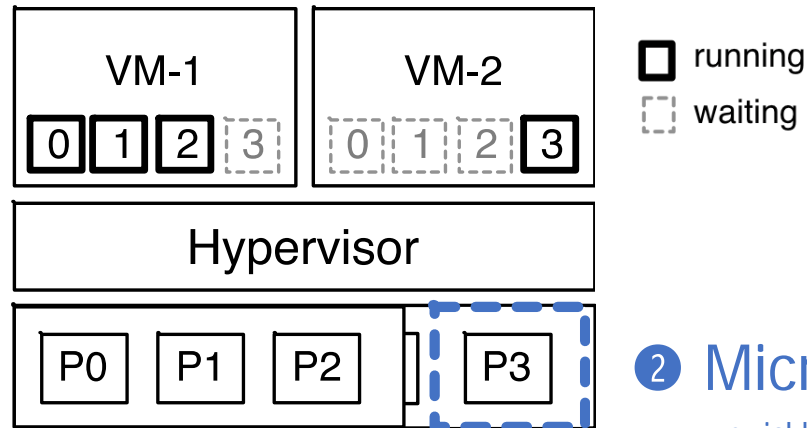


① Normal pool

② Micro-sliced pool
- quickly but briefly schedule vCPUs

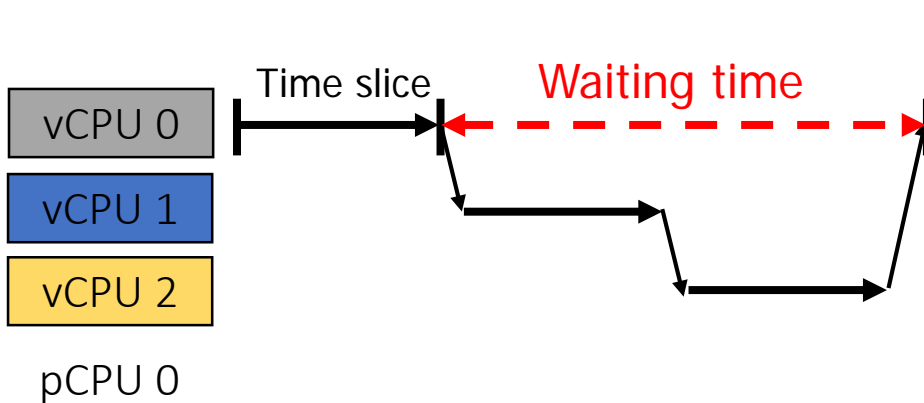


Approach: Dividing CPUs into Two Pools

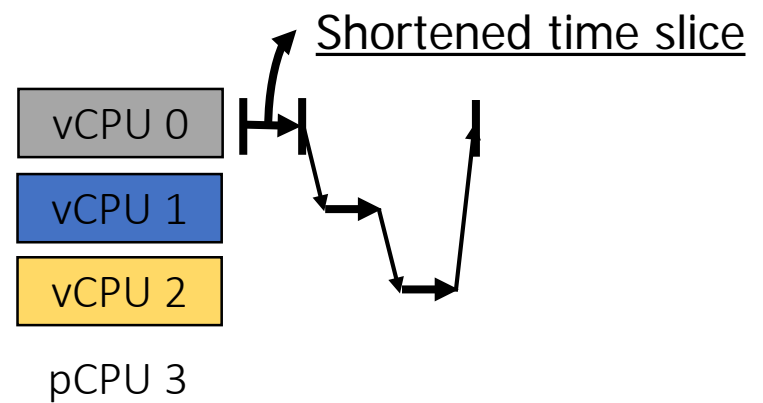


① Normal pool

② Micro-sliced pool
- quickly but briefly schedule vCPUs

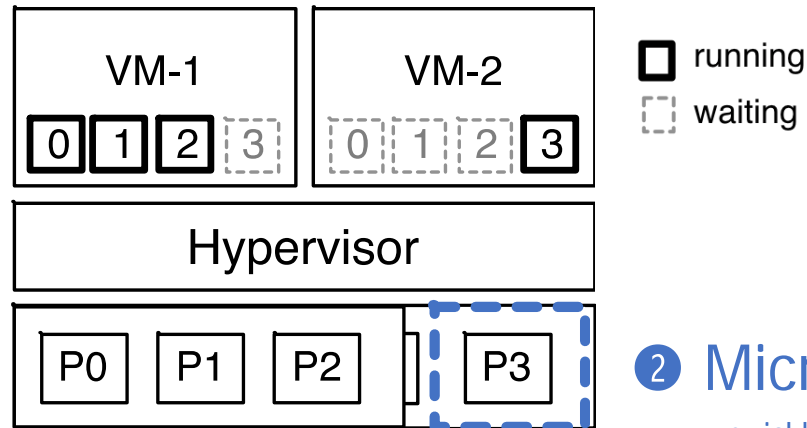


Serving the main work of applications



Serving critical OS services to minimize the *waiting time*

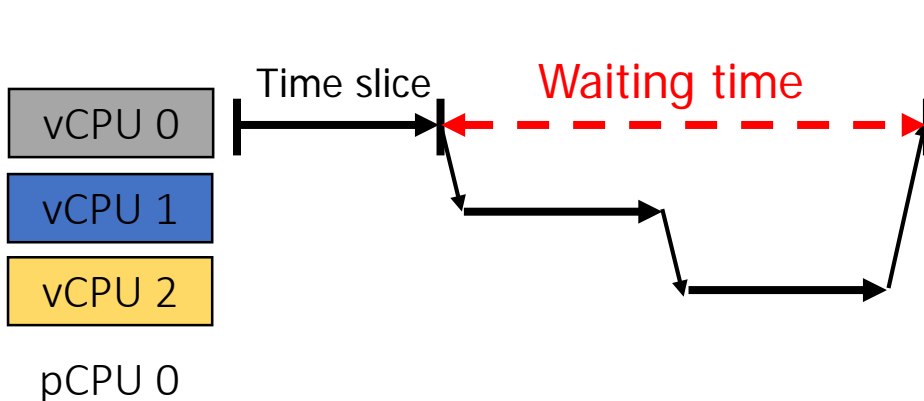
Approach: Dividing CPUs into Two Pools



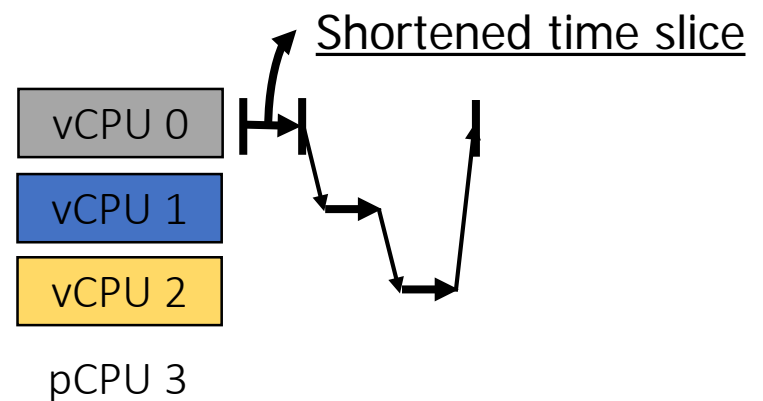
① Normal pool

② Micro-sliced pool

- quickly but briefly schedule vCPUs

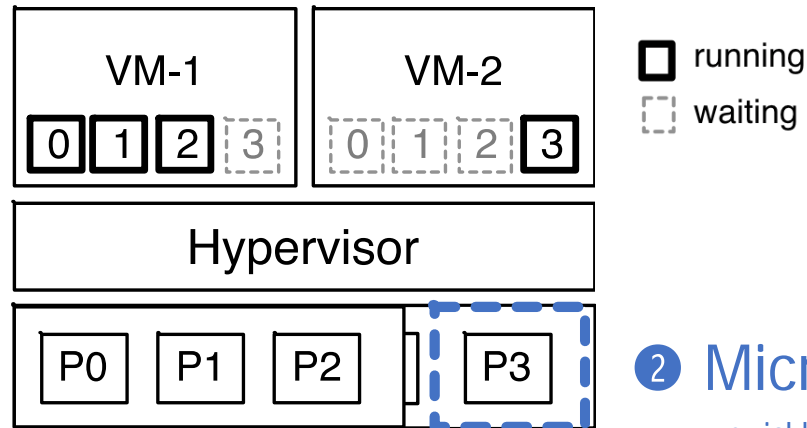


Serving the main work of applications



Serving critical OS services to minimize the *waiting time*

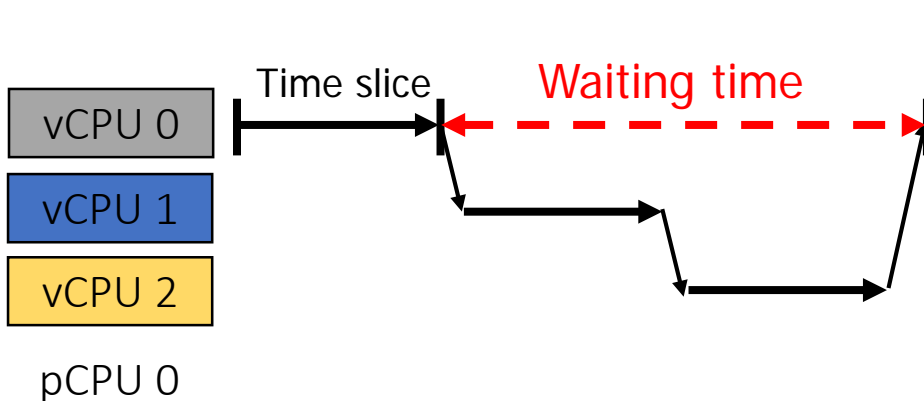
Approach: Dividing CPUs into Two Pools



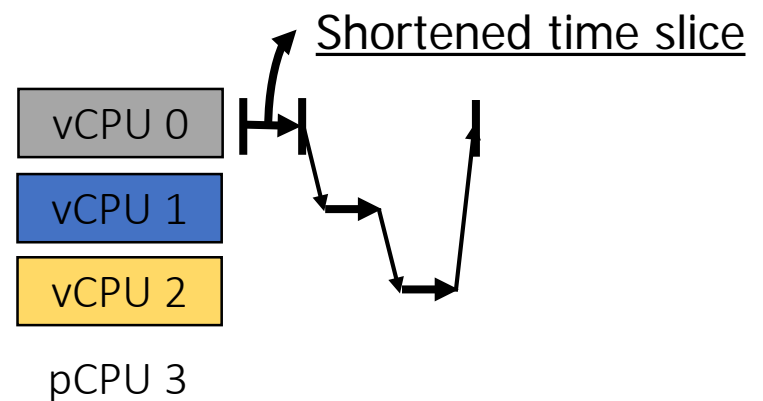
① Normal pool

② Micro-sliced pool

- quickly but briefly schedule vCPUs



Serving the main work of applications



Serving critical OS services to minimize the *waiting time*

Challenges in Serving Critical OS Services on Micro-sliced Cores

Challenges in Serving Critical OS Services on Micro-sliced Cores

1. Precise detection of urgent tasks

Challenges in Serving Critical OS Services on Micro-sliced Cores

1. Precise detection of urgent tasks

2. Guest OS transparency

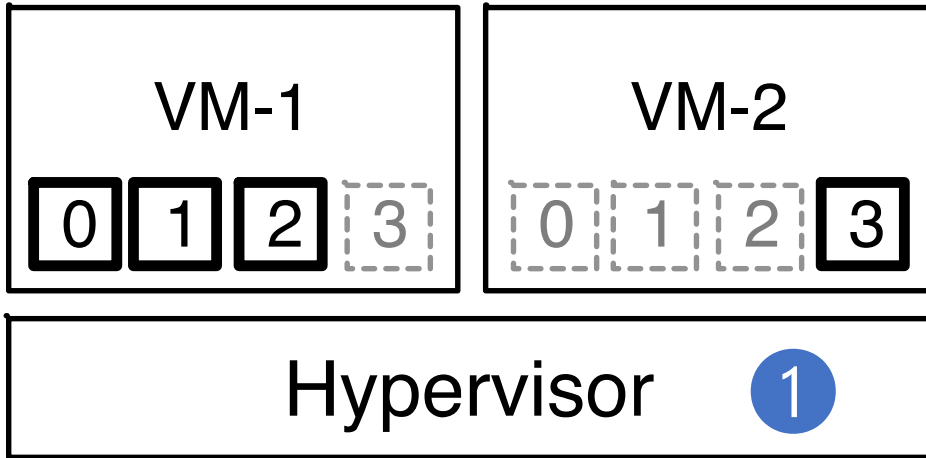
Challenges in Serving Critical OS Services on Micro-sliced Cores

1. Precise detection of urgent tasks

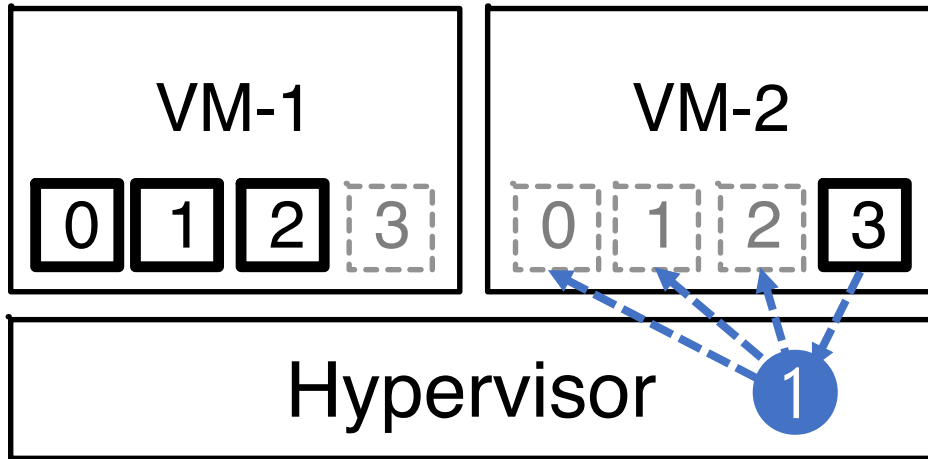
2. Guest OS transparency

3. Dynamic adjustment of micro-sliced cores

Detecting Critical OS Services



Detecting Critical OS Services



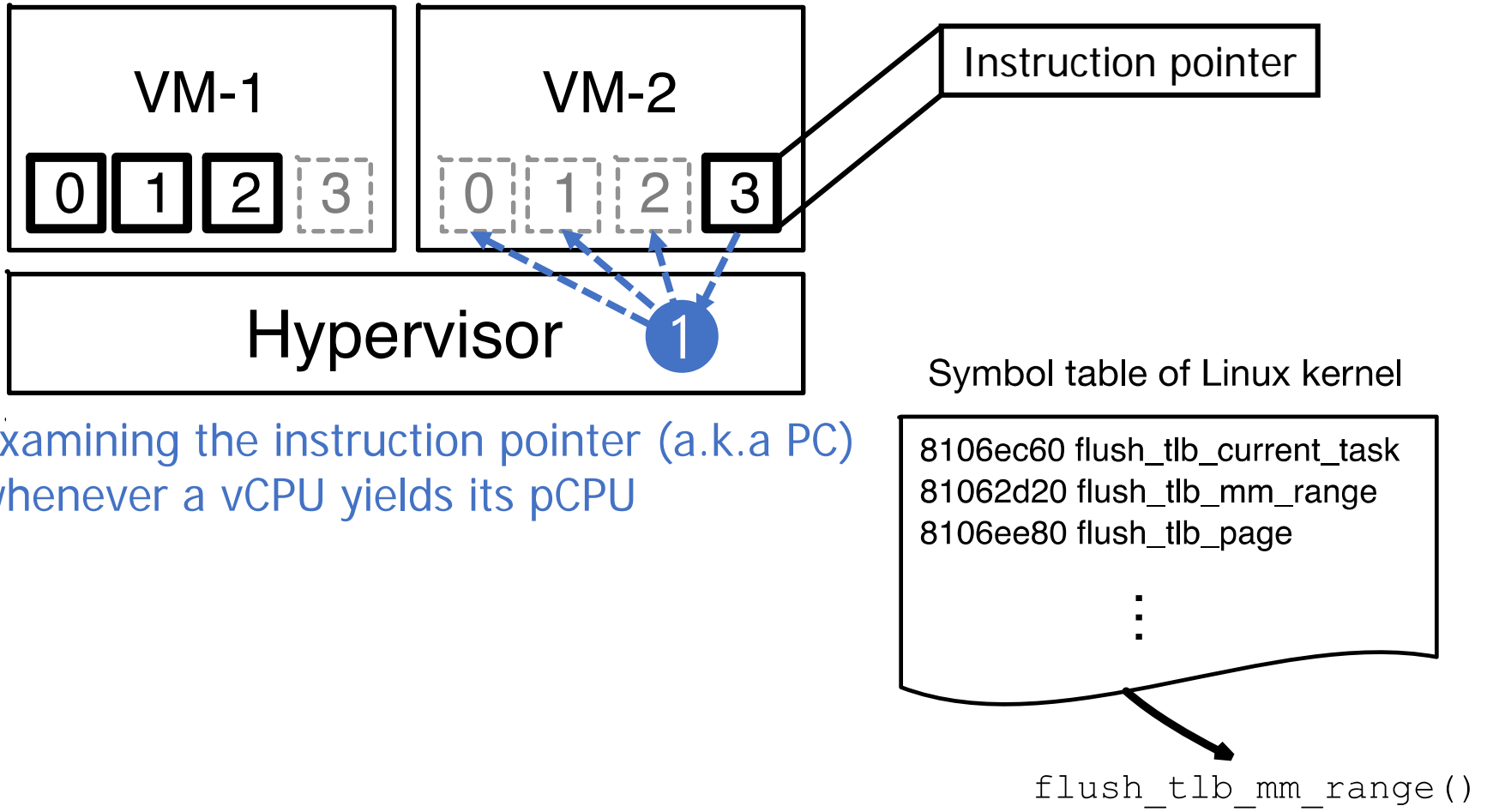
Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

Symbol table of Linux kernel

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

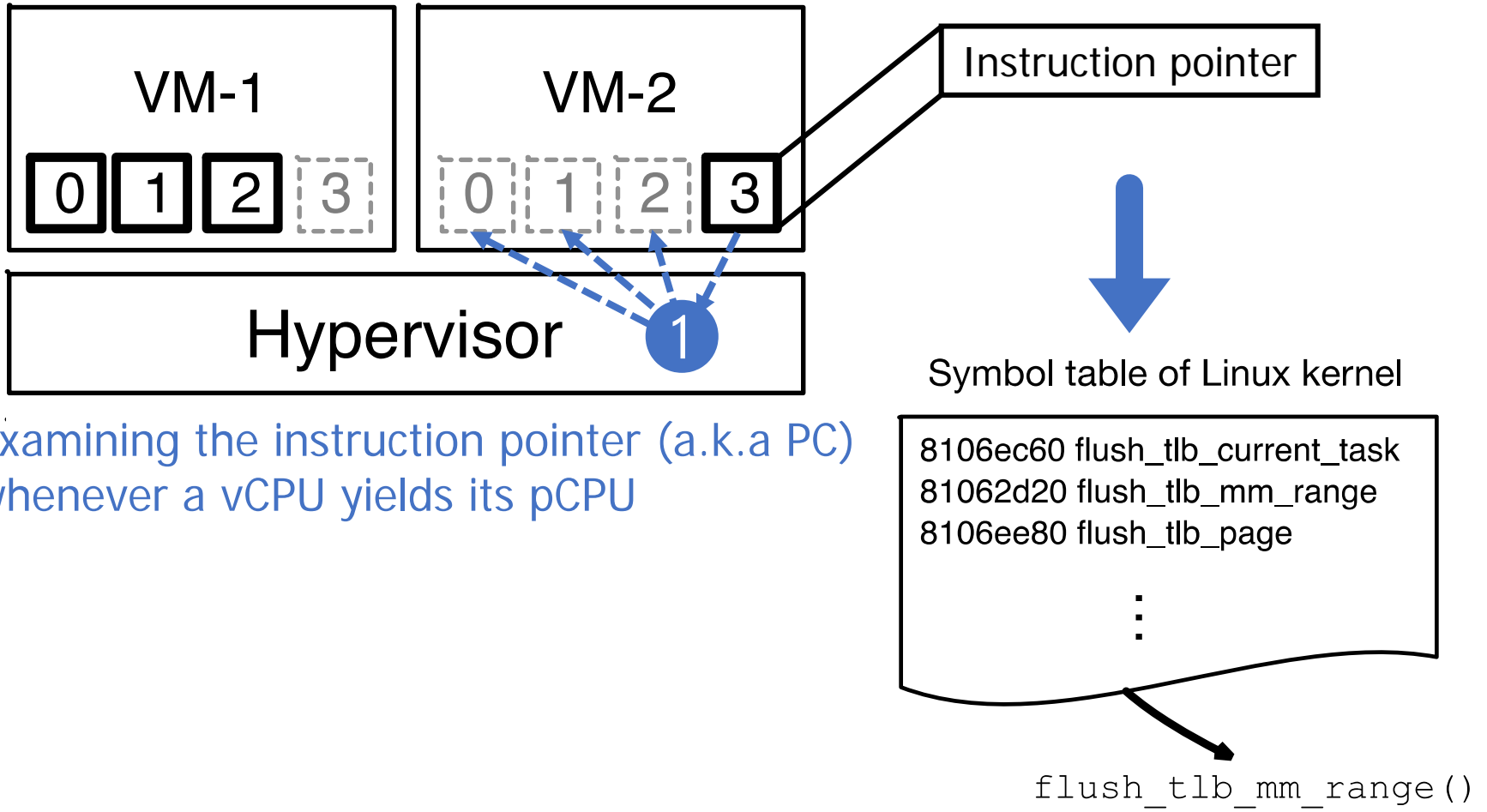
flush_tlb_mm_range()

Detecting Critical OS Services



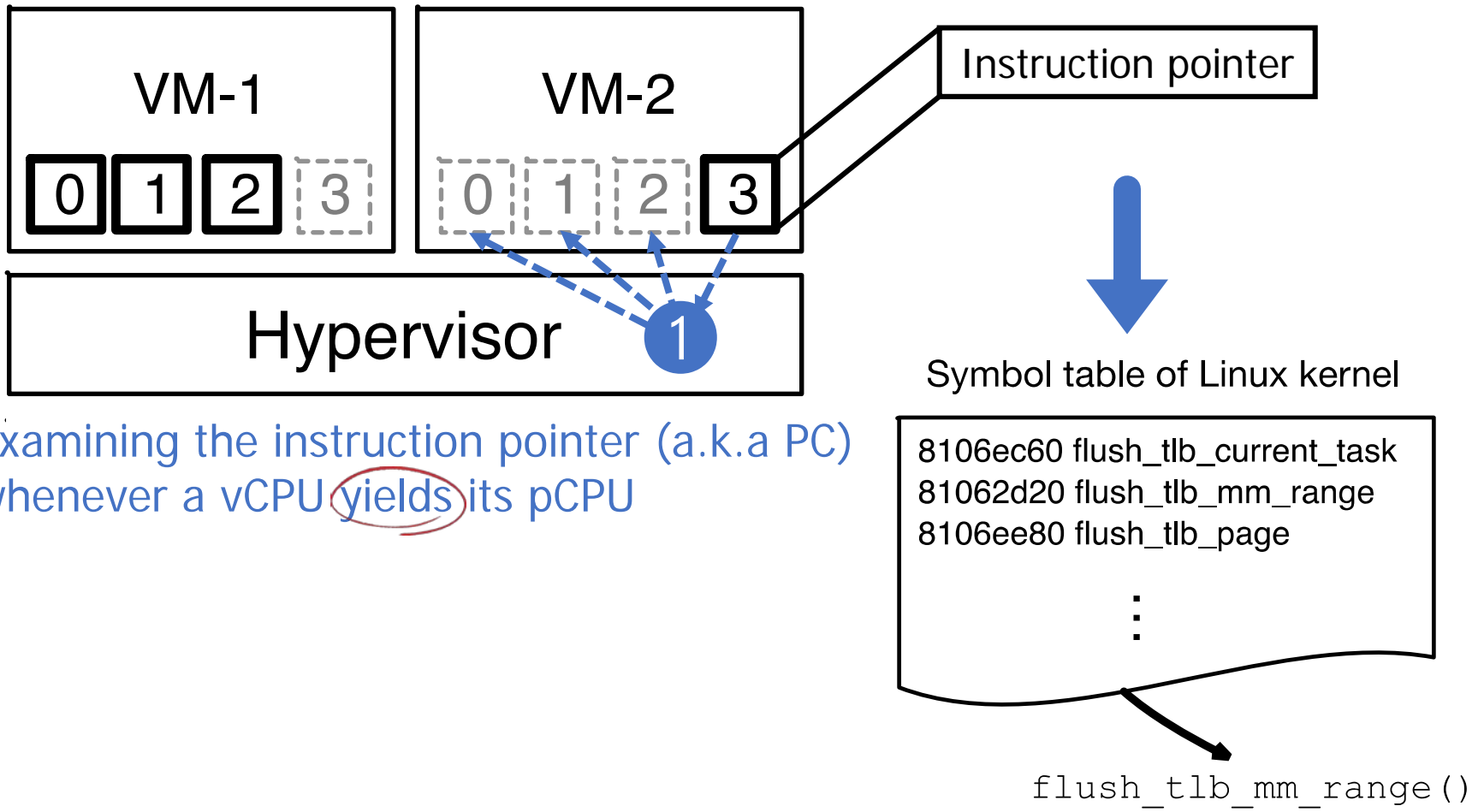
Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

Detecting Critical OS Services

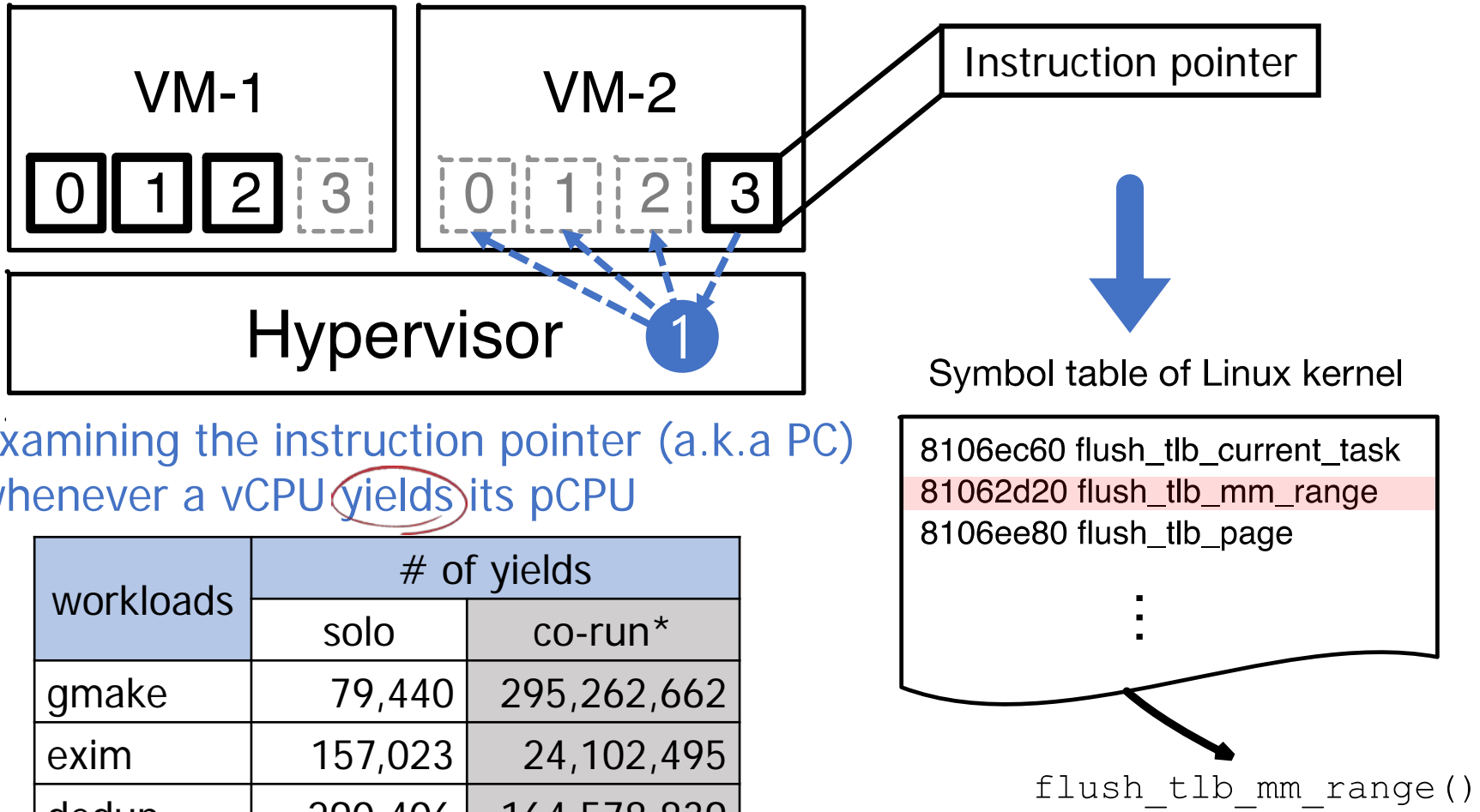


Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

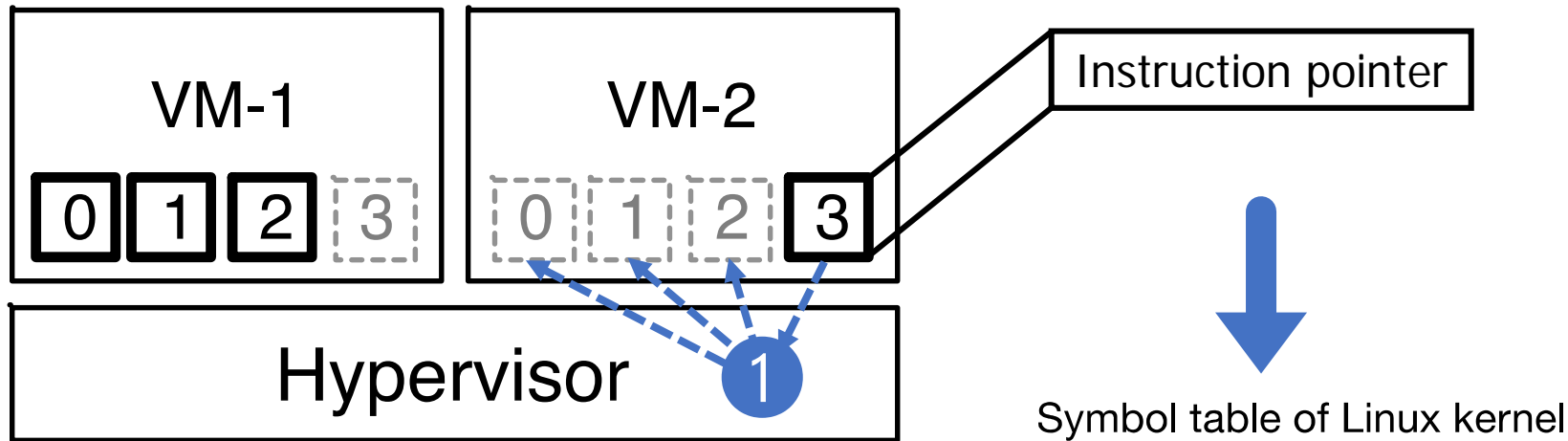
Detecting Critical OS Services



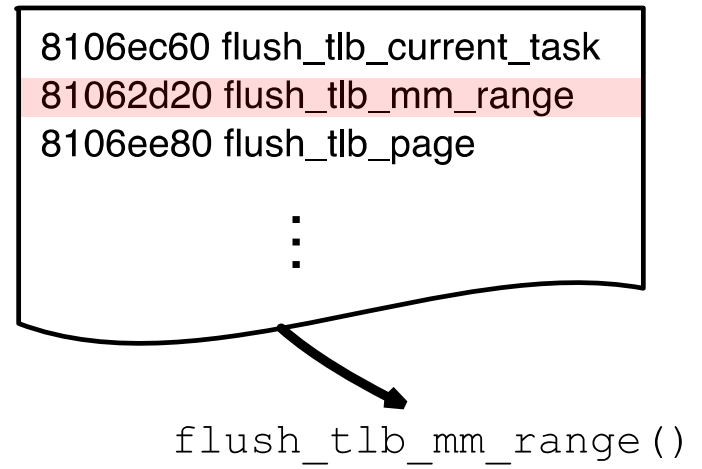
Detecting Critical OS Services



Detecting Critical OS Services



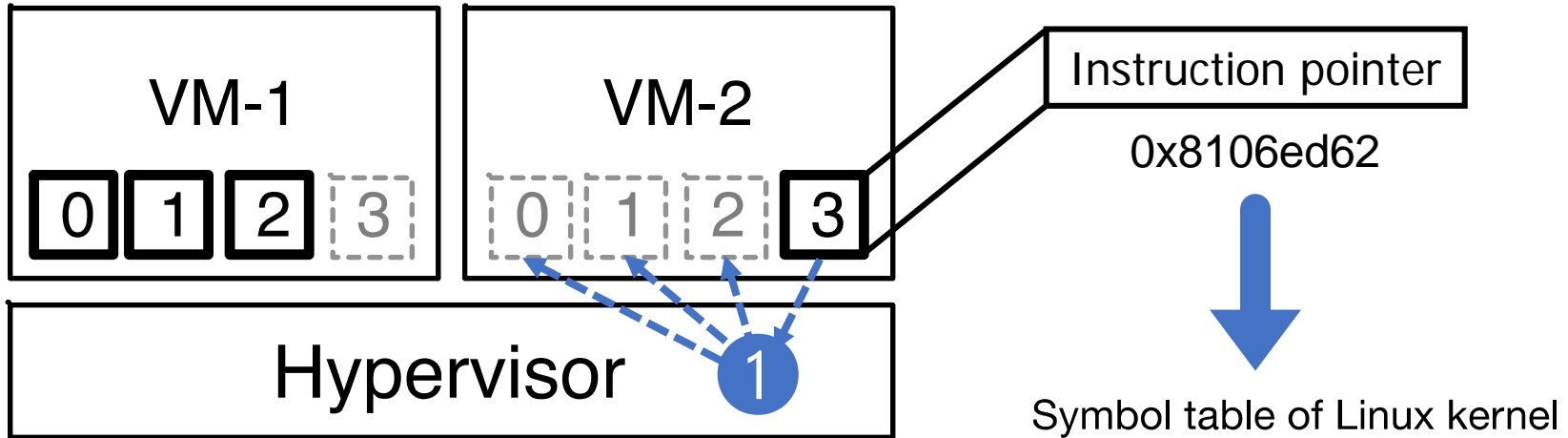
Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU



workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

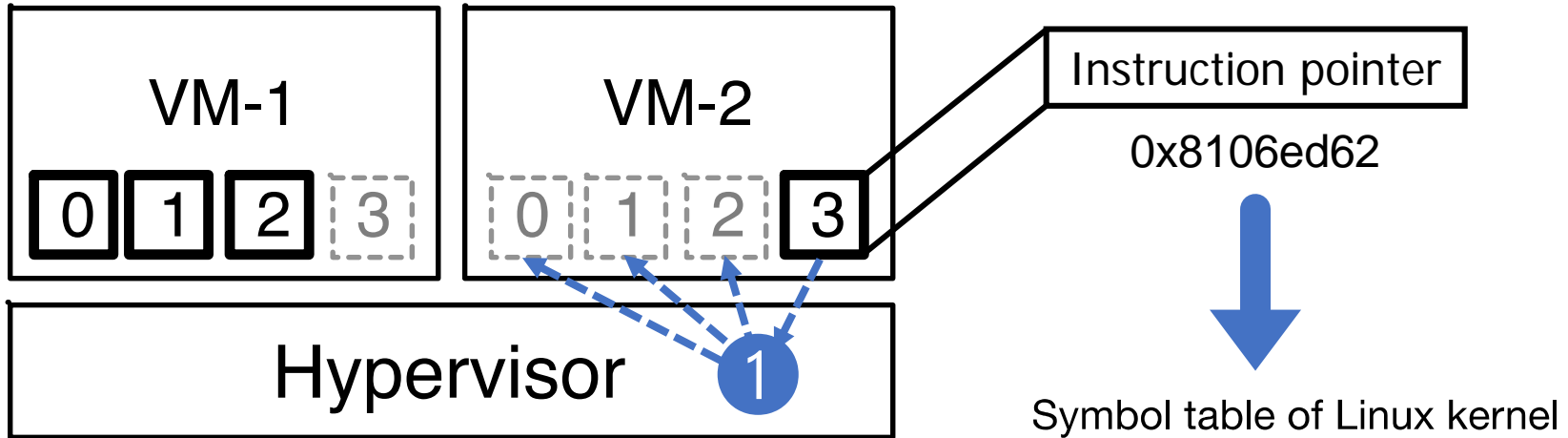
flush_tlb_mm_range()

May need to inspect vCPU siblings

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
⋮	

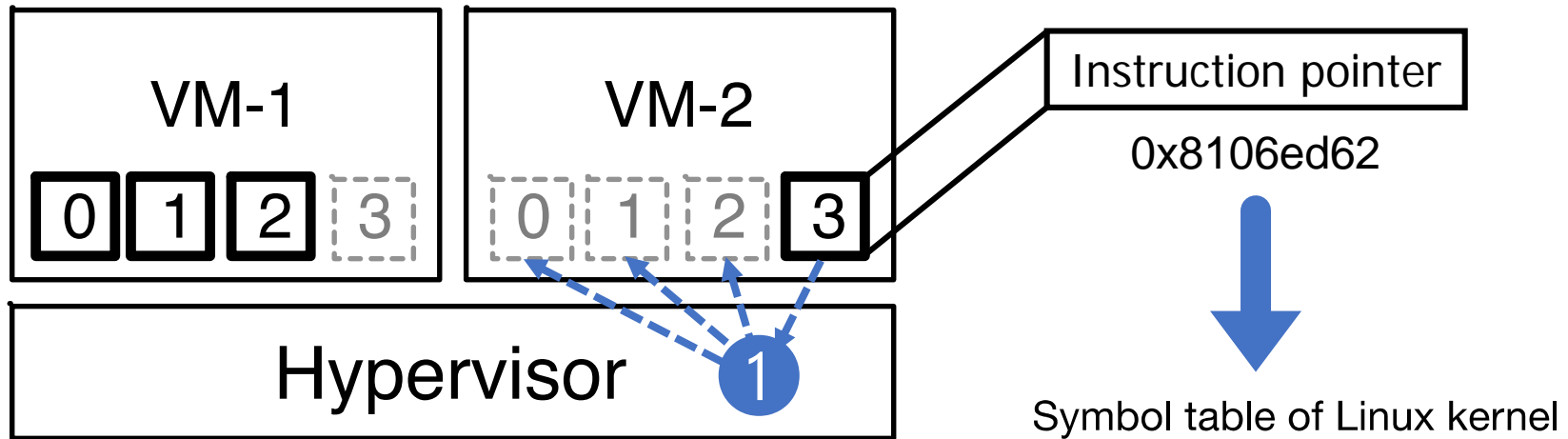
flush_tlb_mm_range()

May need to inspect vCPU siblings

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

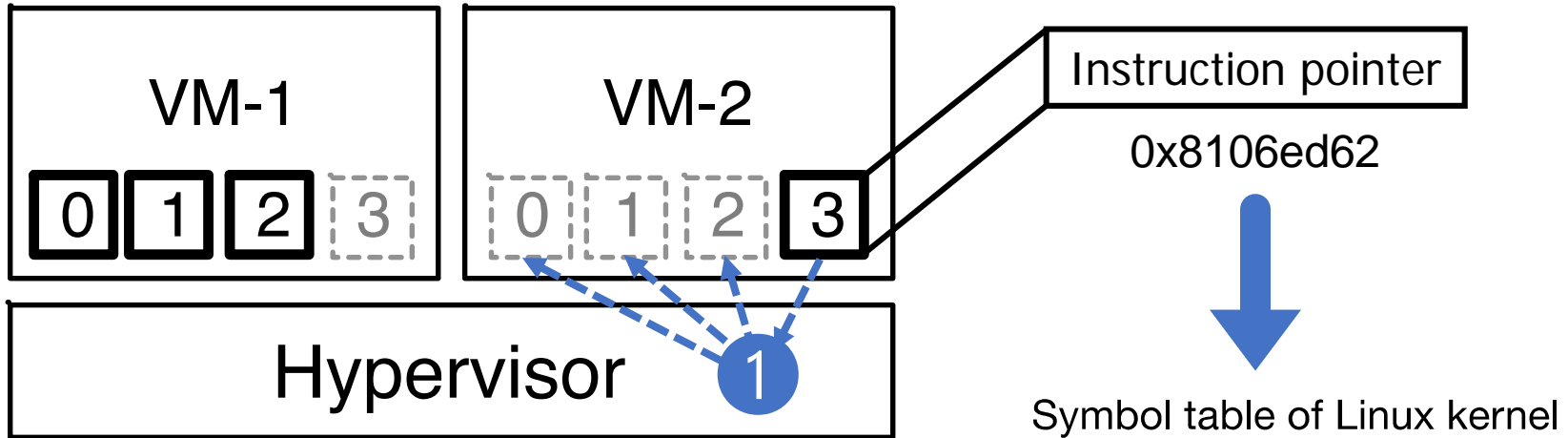
flush_tlb_mm_range()

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

May need to inspect vCPU siblings

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

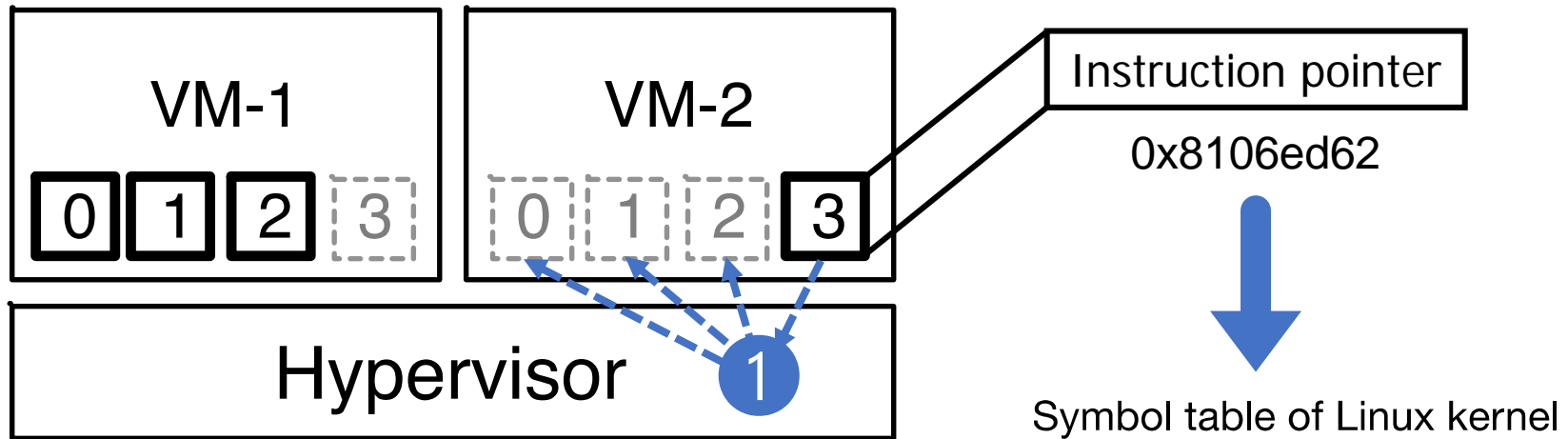
flush_tlb_mm_range()

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

May need to inspect vCPU siblings

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

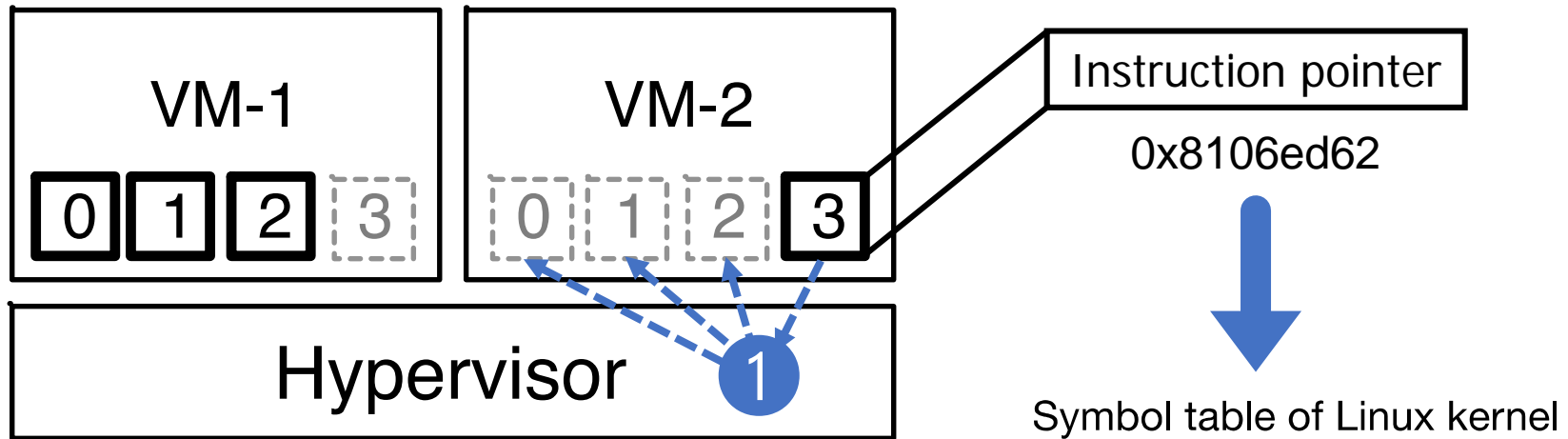
flush_tlb_mm_range()

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

May need to inspect vCPU siblings

Detecting Critical OS Services



Examining the instruction pointer (a.k.a PC) whenever a vCPU yields its pCPU

8106ec60	flush_tlb_current_task
81062d20	flush_tlb_mm_range
8106ee80	flush_tlb_page
	⋮

flush_tlb_mm_range()

May need to inspect vCPU siblings

workloads	# of yields	
	solo	co-run*
gmake	79,440	295,262,662
exim	157,023	24,102,495
dedup	290,406	164,578,839
vips	644,643	57,650,538

* Concurrently running with Swaptions of PARSEC

Profiling Virtual CPU Scheduling Logs

- Investigating frequently preempted regions



vCPU scheduling trace
(w/ Inst. Pointer)

```
System.map-4.4.0-104-generic
ffffffff811045c0 t hotplug_cfd
ffffffff81104680 t generic_exec_single
ffffffff811047a0 T smp_call_function_single
ffffffff811048d0 T smp_call_function_single_async
ffffffff81104940 T smp_call_function_any
ffffffff81104a10 T smp_call_function_many
ffffffff81104c70 T smp_call_function
ffffffff81104ca0 T kick_all_cpus_sync
ffffffff81104cd0 T wake_up_all_idle_cpus
ffffffff81104d20 T on_each_cpu
ffffffff81104d80 T on_each_cpu_mask
ffffffff81104de0 T on_each_cpu_cond
<em.map-4.4.0-104-generic CWD: /boot Line: 6108
```

Kernel symbol tables

Critical Guest OS Components

Module	Operation
sched	scheduler_ipi()
	resched_curr() ...
mm	flush_tlb_all()
	get_page_from_freelist() ...
irq	irq_enter()
	irq_exit() ...
spinlock	__raw_spin_unlock()
	__raw_spin_unlock_irq() ...

In our paper, you can find the table in details

Profiling Virtual CPU Scheduling Logs

- Investigating frequently preempted regions



vCPU scheduling trace
(w/ Inst. Pointer)

Critical Guest OS Components

Module	Operation
	scheduler_ipi()
sched	resched_curr()

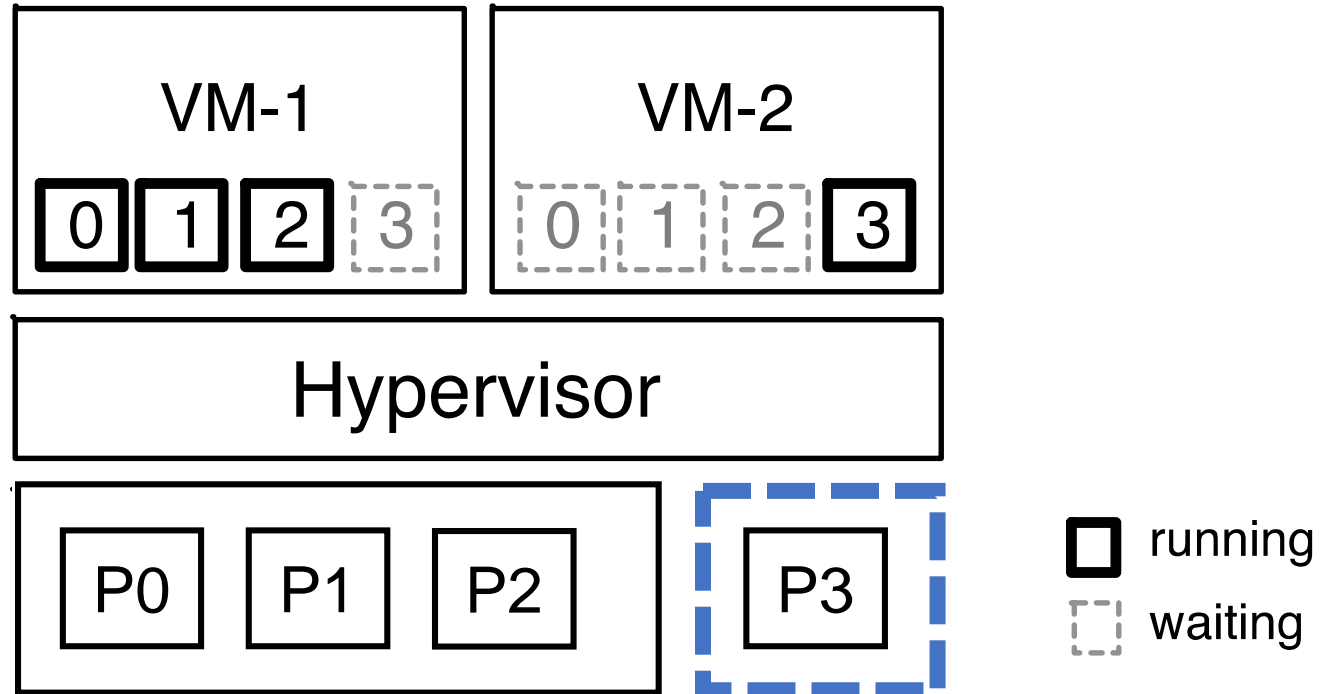
Instruction pointer and kernel symbols enable to precisely detect vCPUs preempted while executing critical OS services without guest OS modification

SPINLOCK	__raw_spin_unlock_irq()
	...

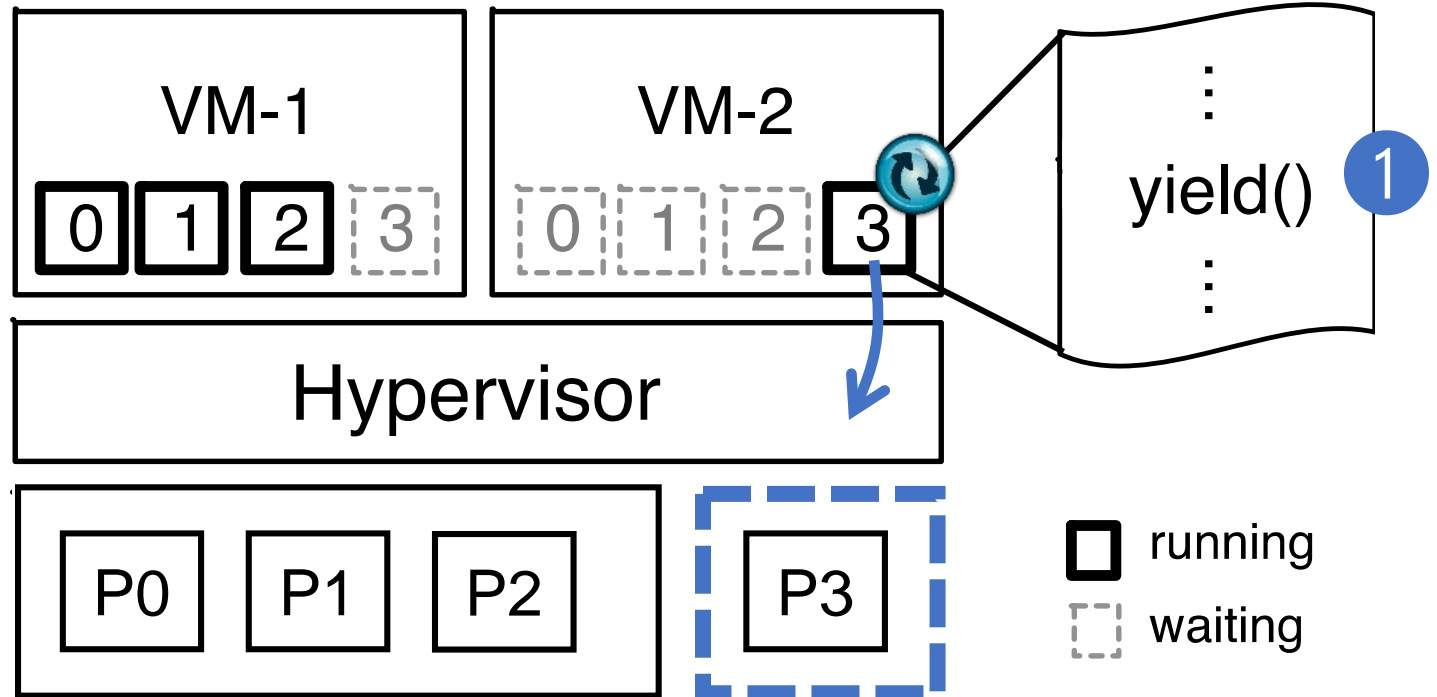
In our paper, you can find the table in details

Kernel symbol tables

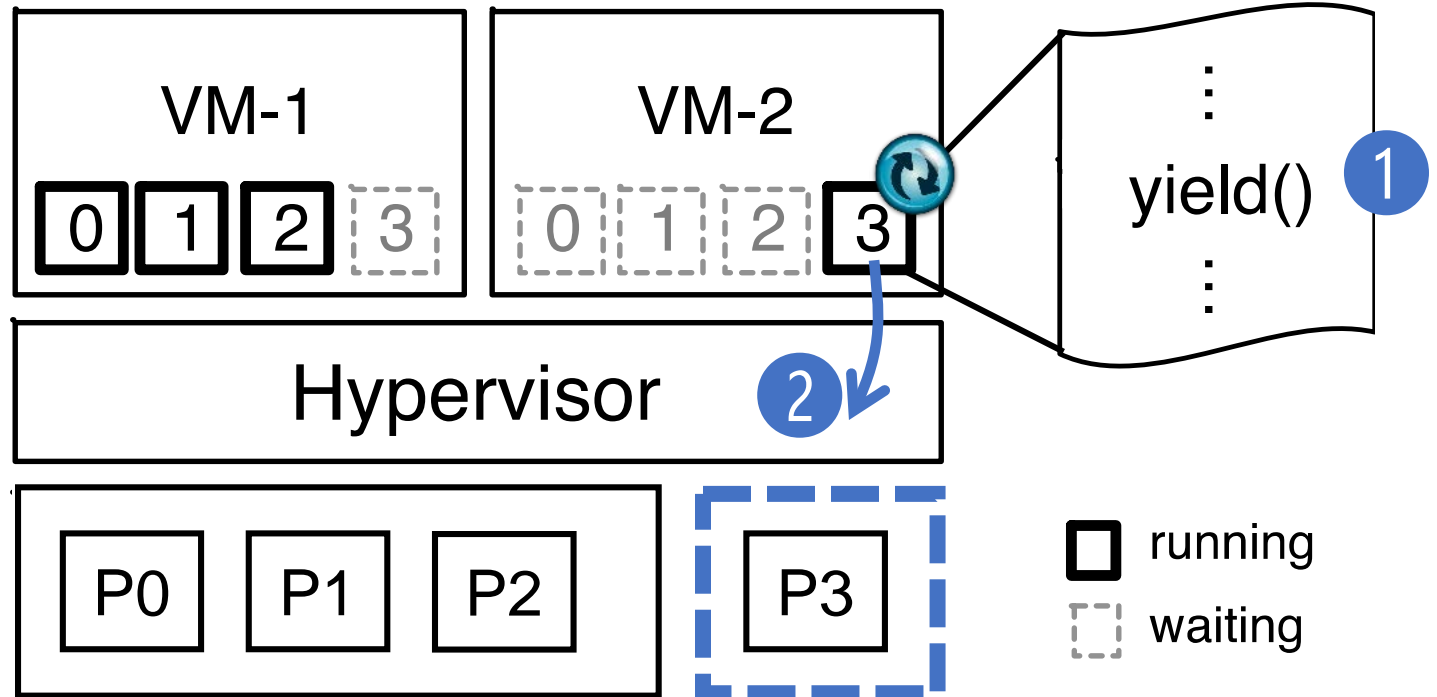
Accelerating Critical Sections



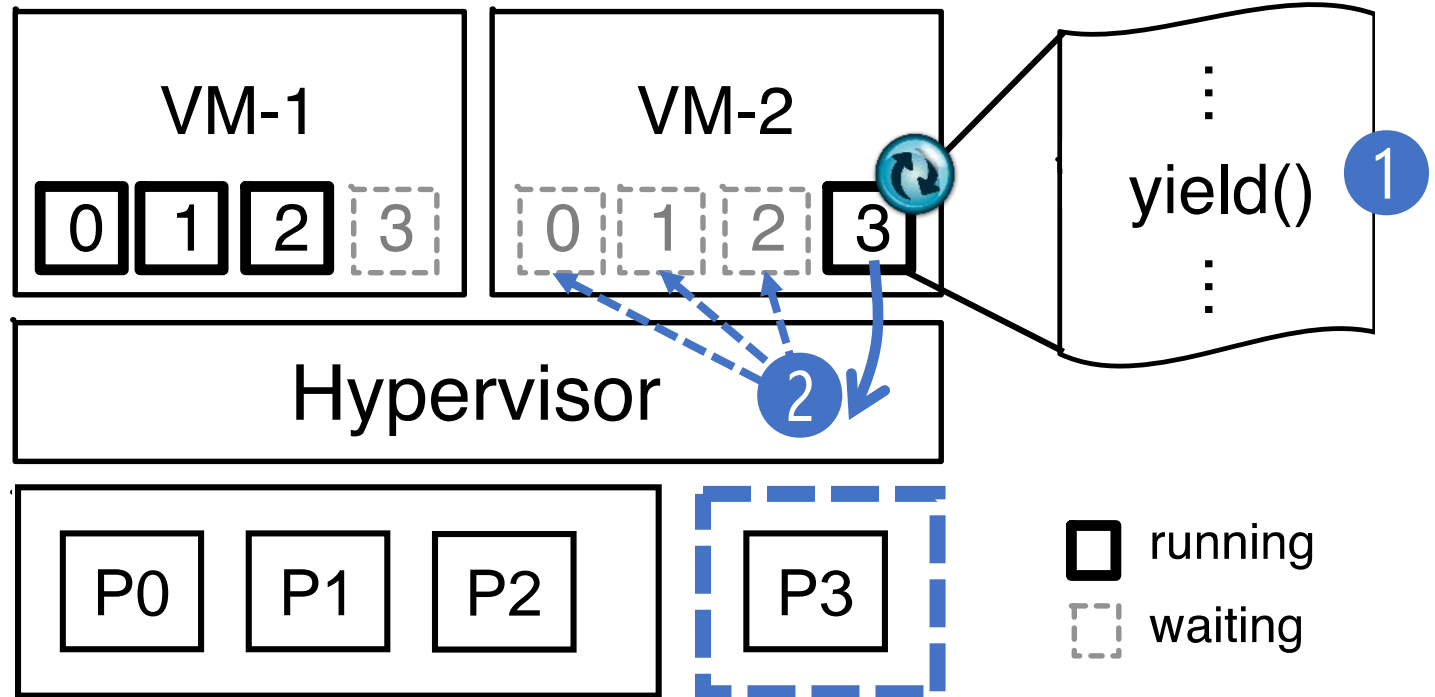
Accelerating Critical Sections



Accelerating Critical Sections

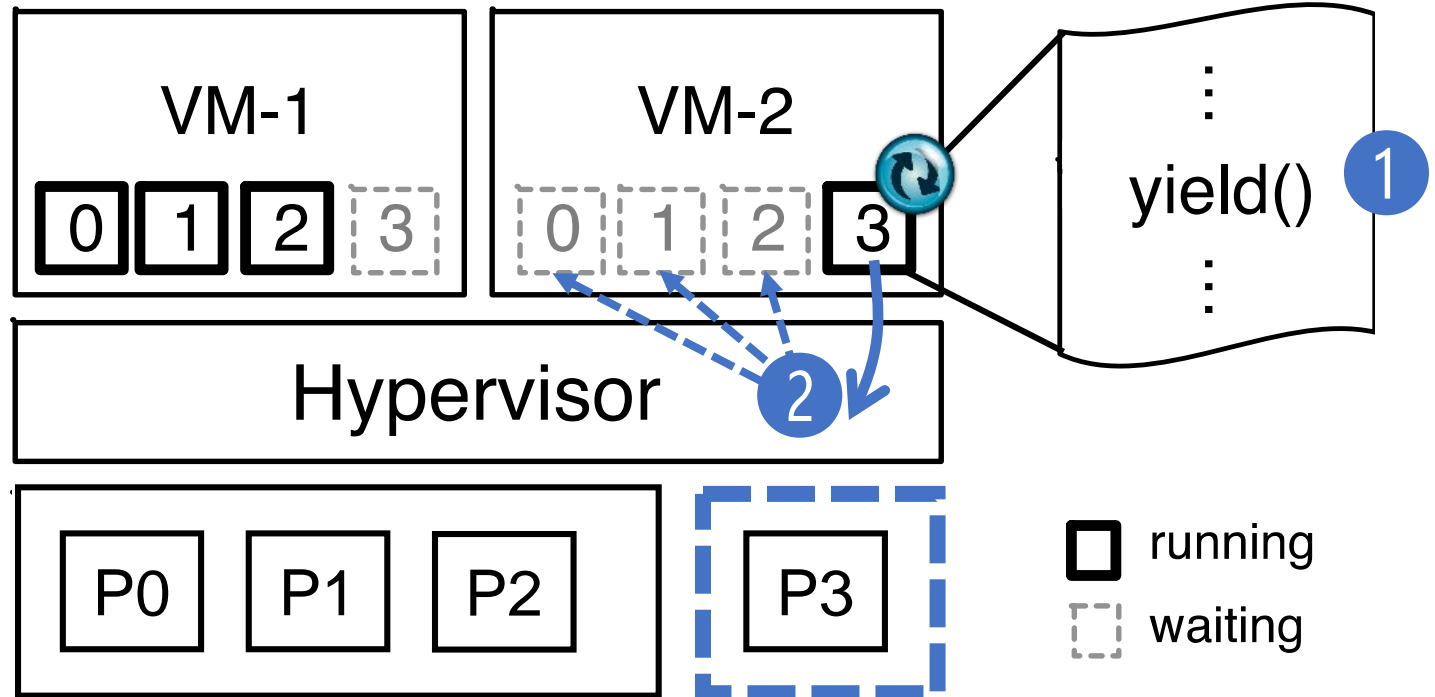


Accelerating Critical Sections



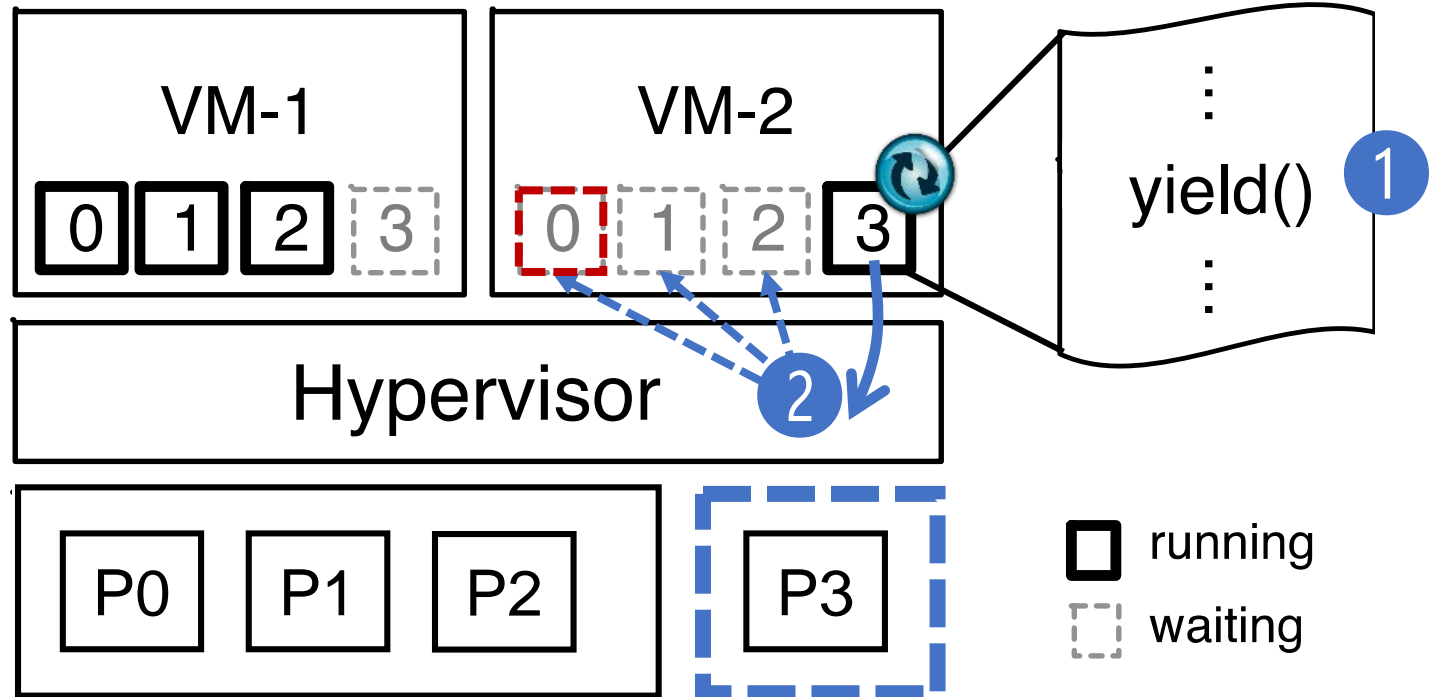
① Yield occurring

Accelerating Critical Sections



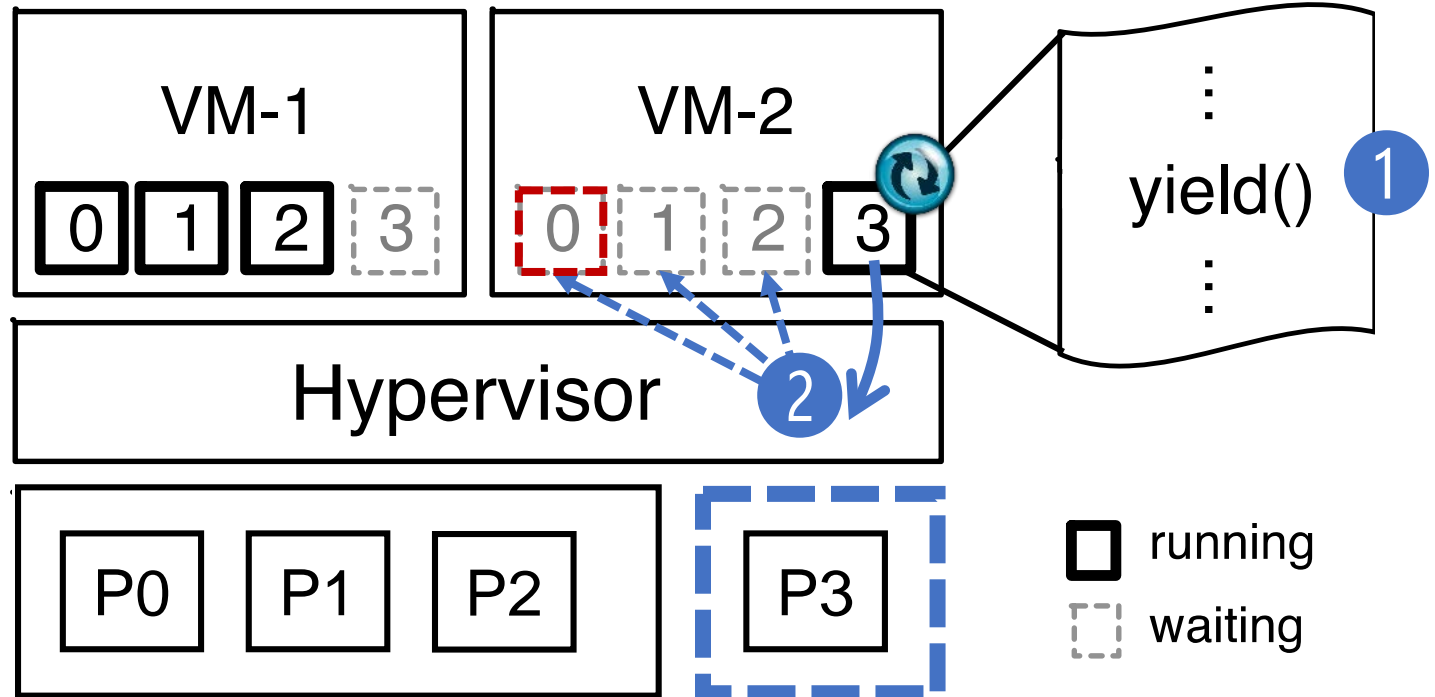
① Yield occurring

Accelerating Critical Sections



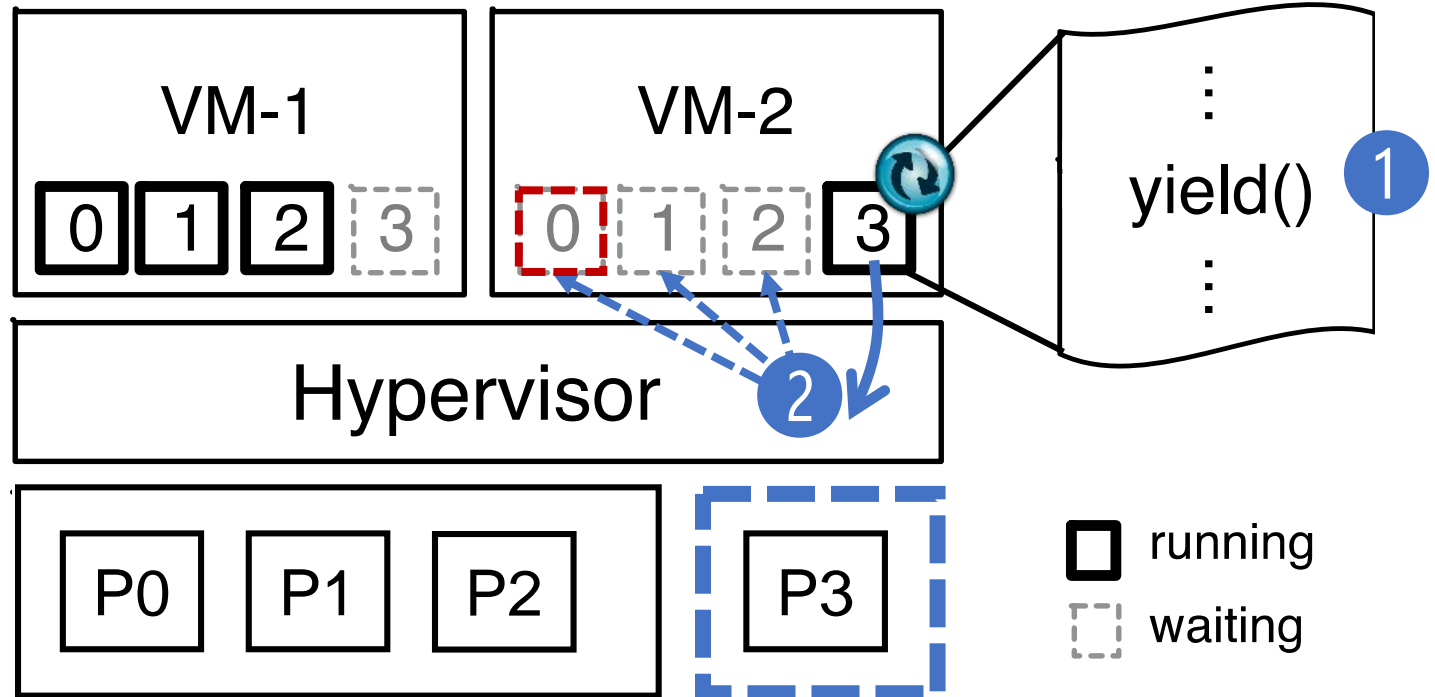
- 1 Yield occurring
- 2 Investigating the preempted vCPUs

Accelerating Critical Sections



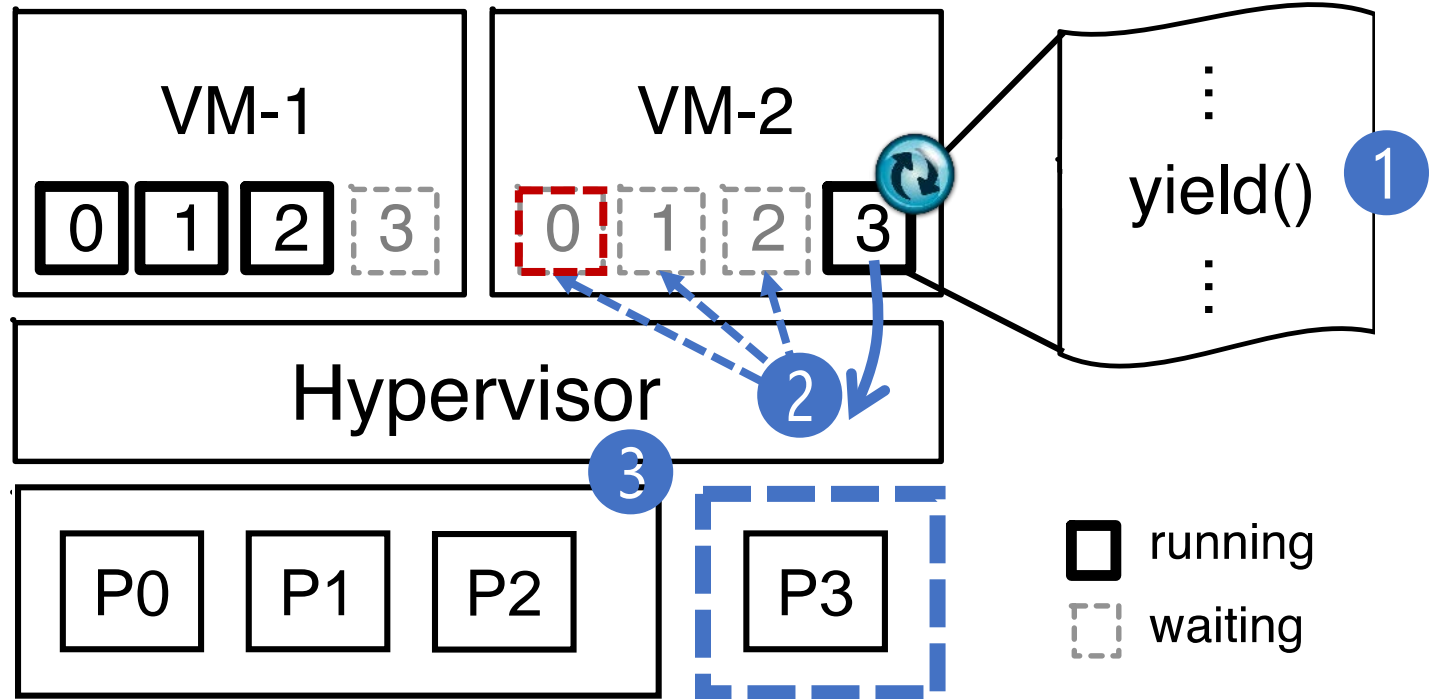
- 1 Yield occurring
- 2 Investigating the preempted vCPUs

Accelerating Critical Sections



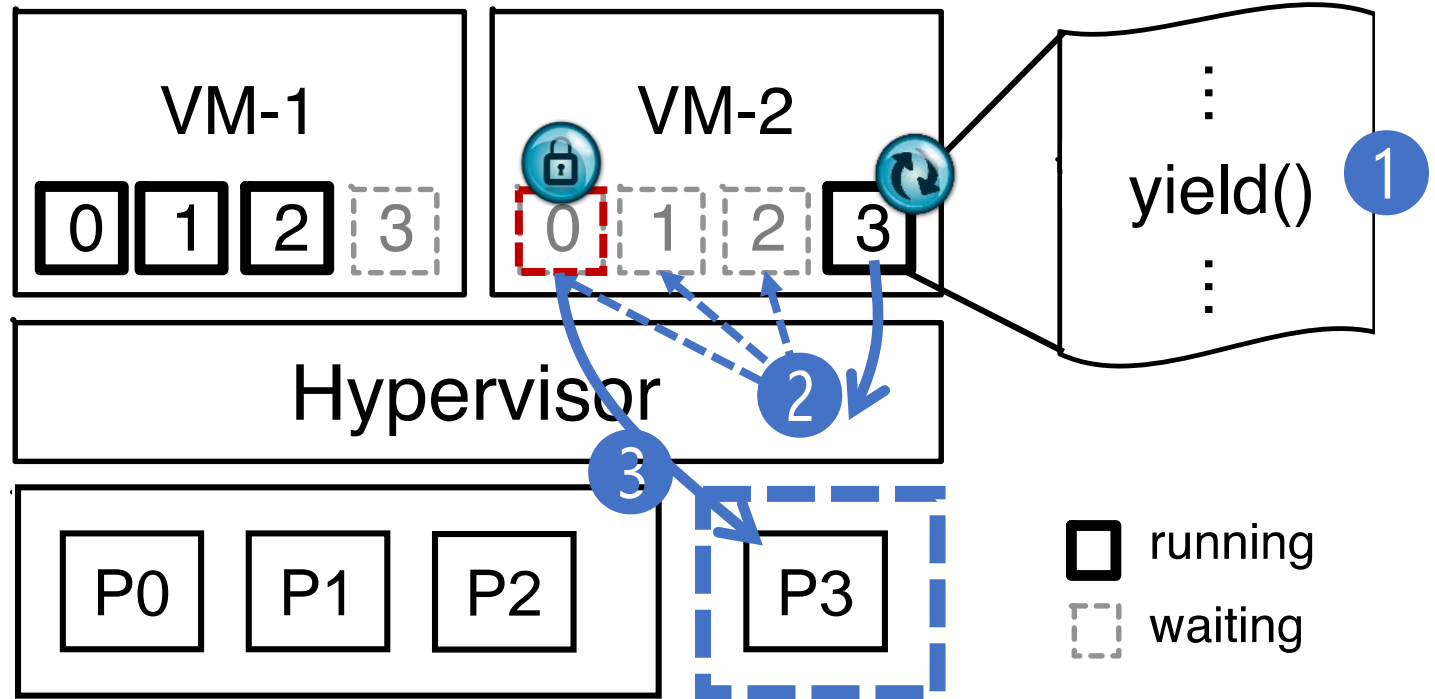
- 1 Yield occurring
- 2 Investigating the preempted vCPUs

Accelerating Critical Sections



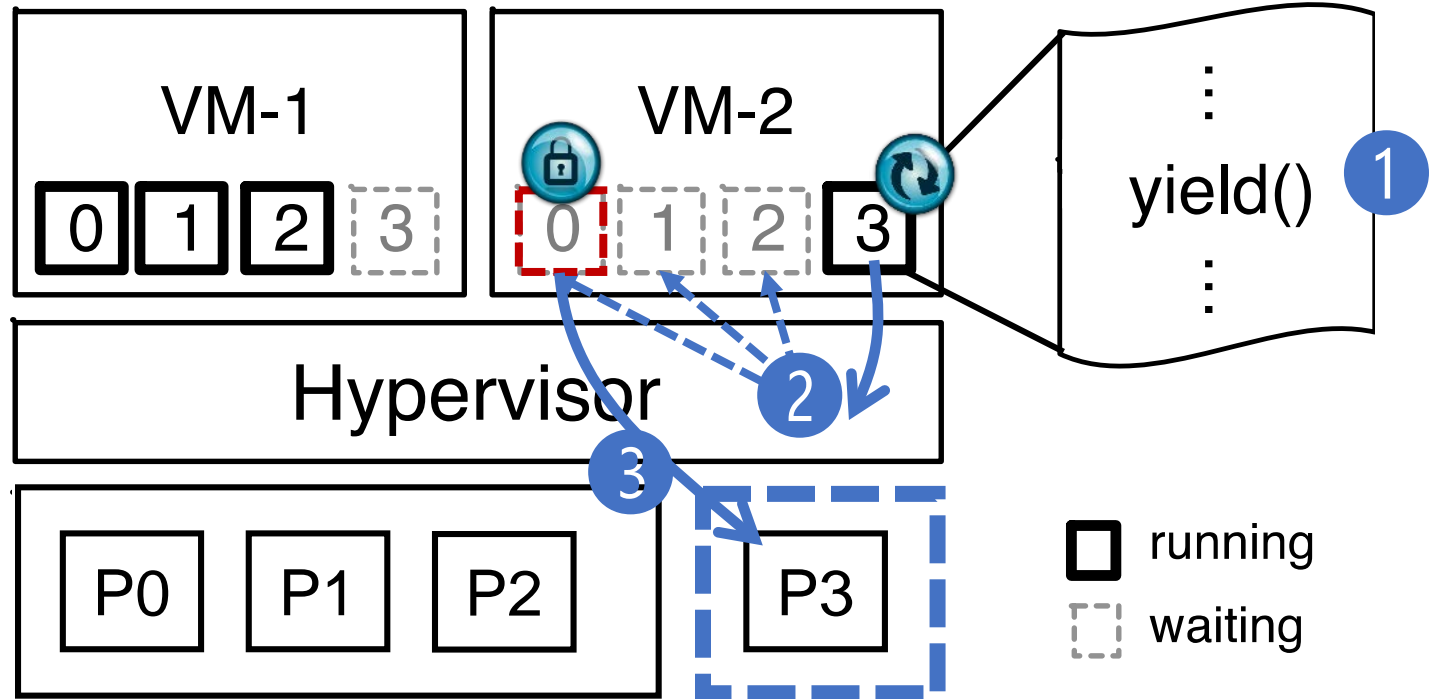
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical Sections



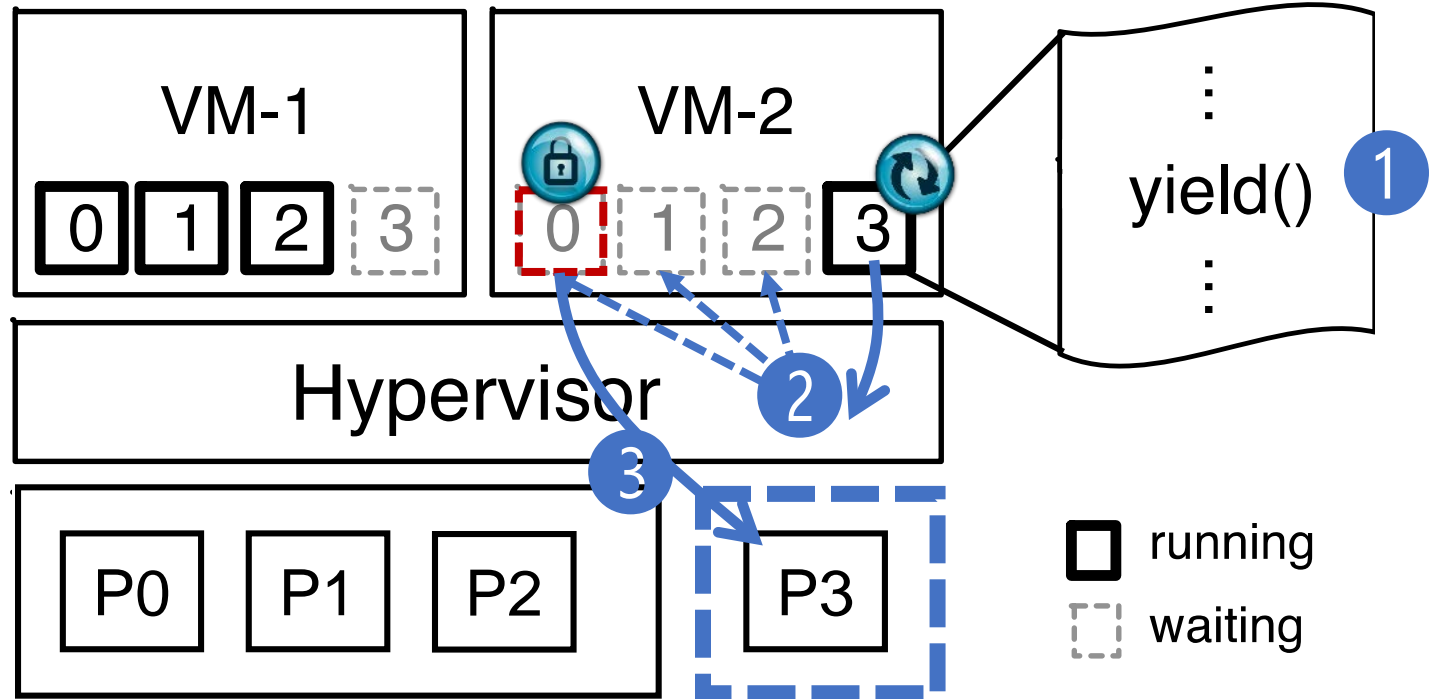
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical Sections



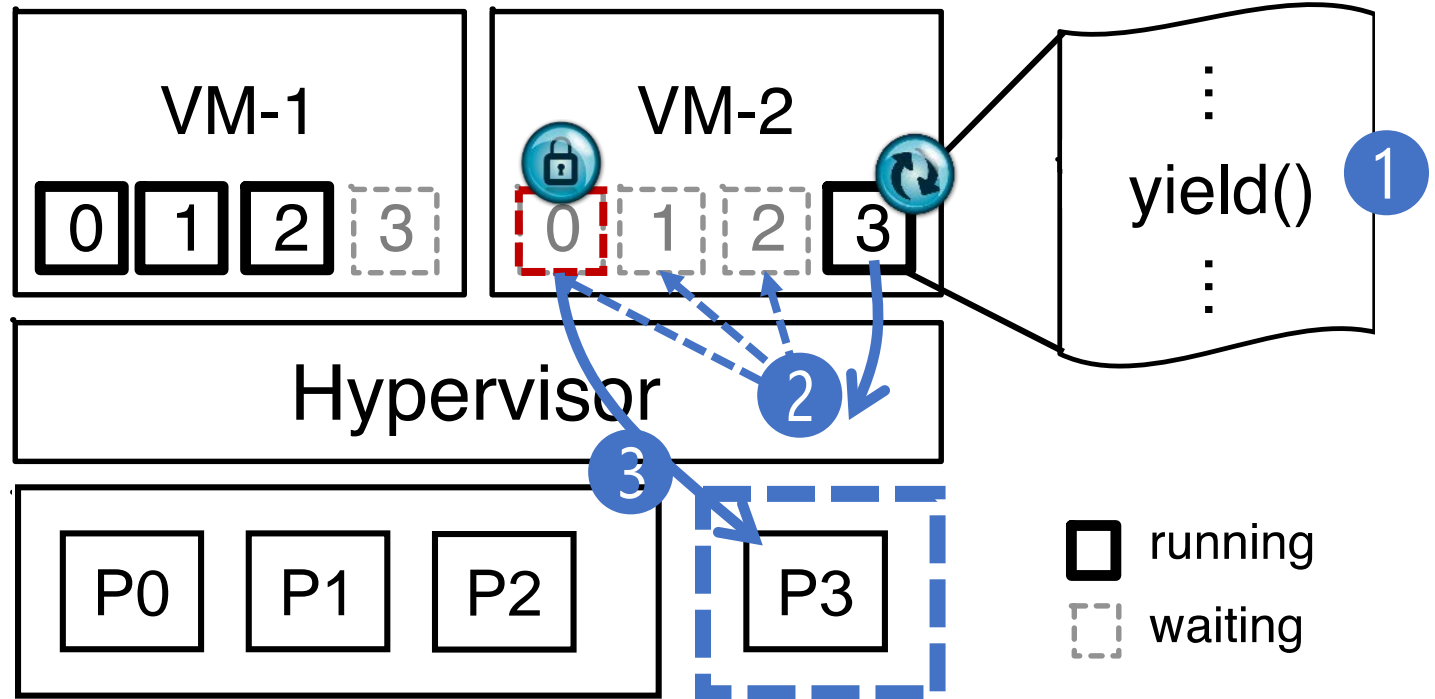
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical Sections



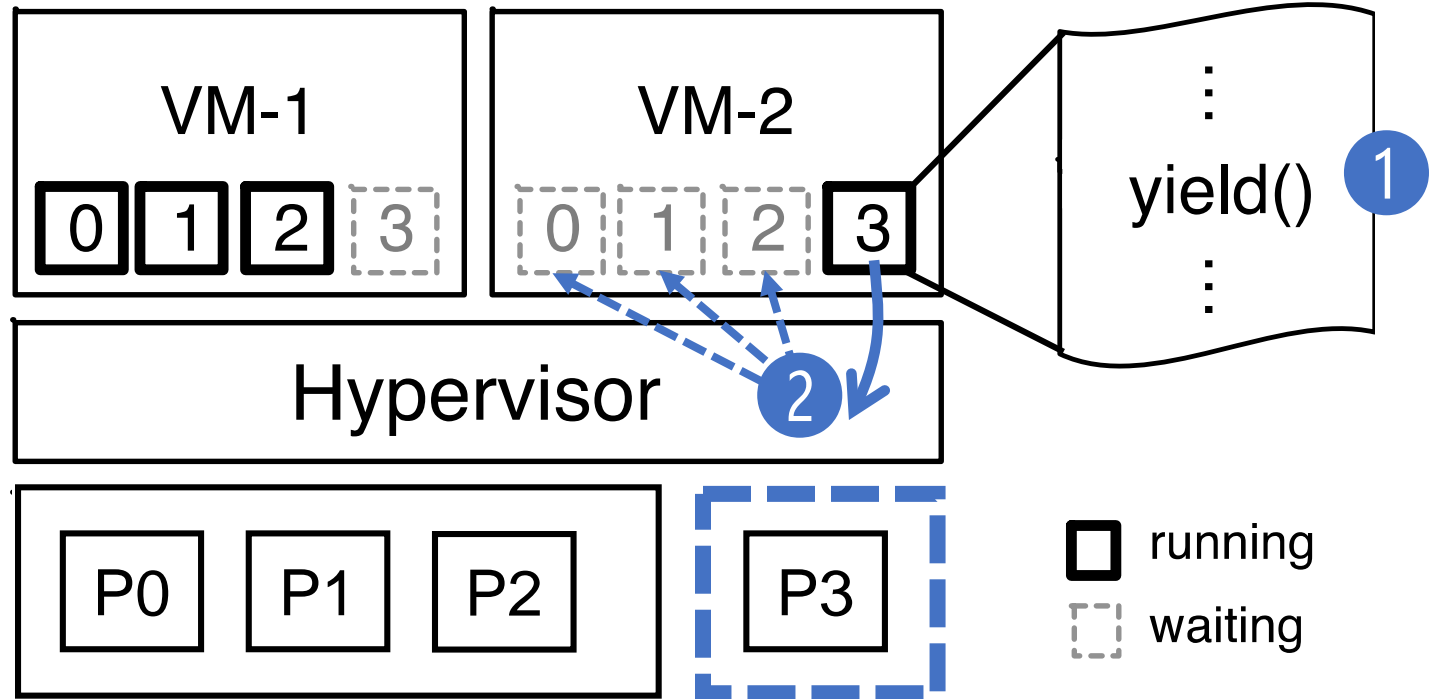
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical Sections



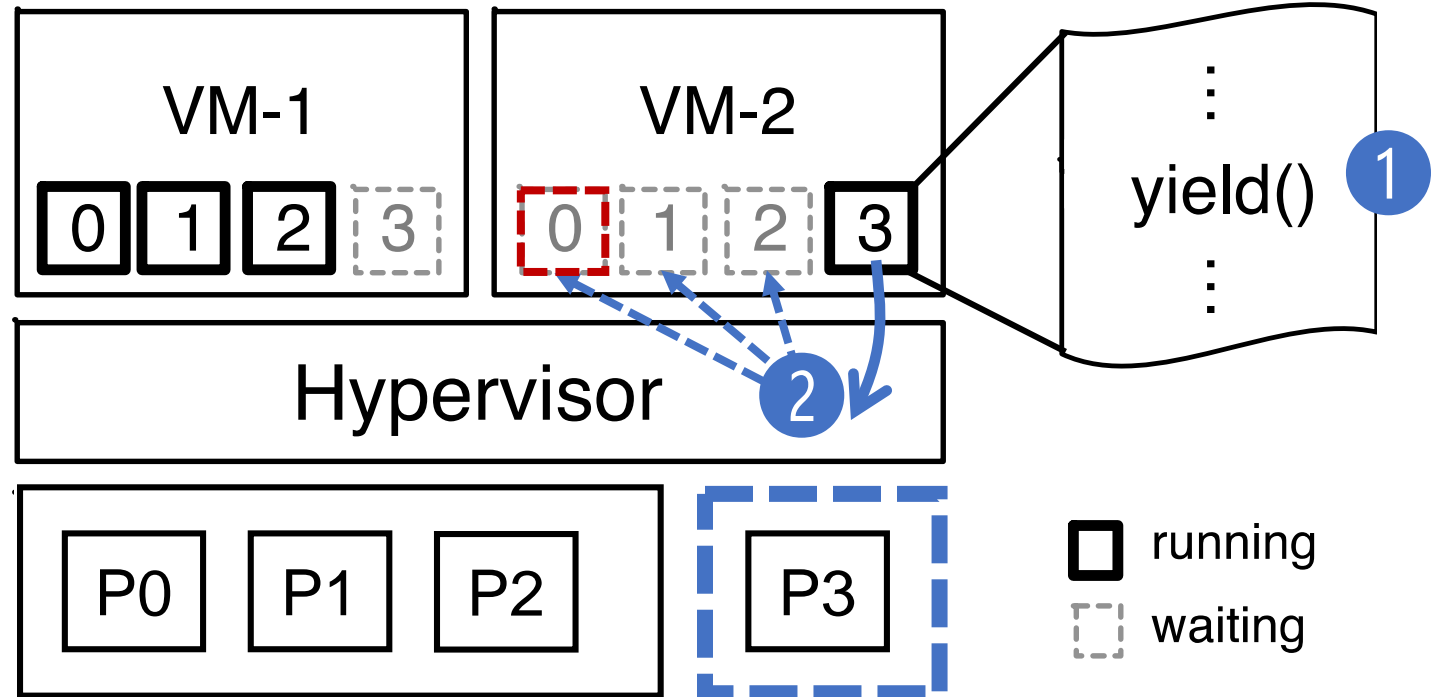
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical TLB Synchronizations



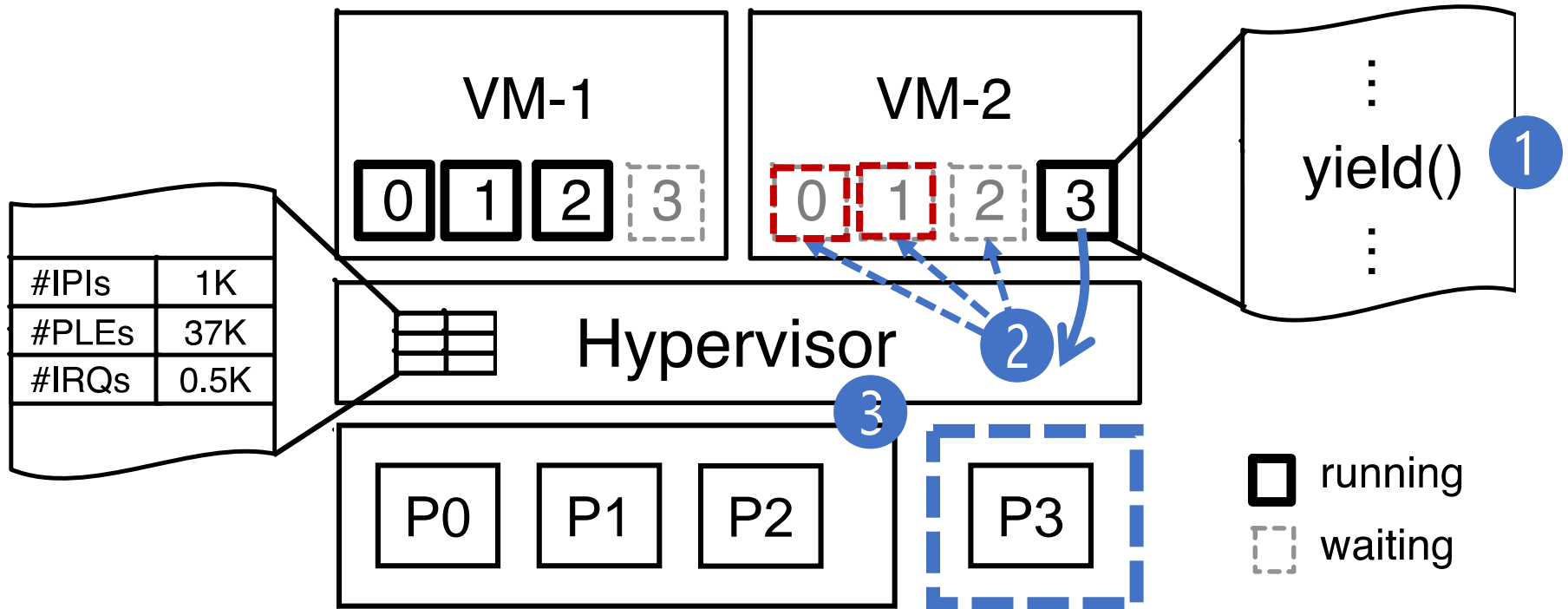
- 1 Yield occurring
- 2 Investigating the preempted vCPUs

Accelerating Critical TLB Synchronizations



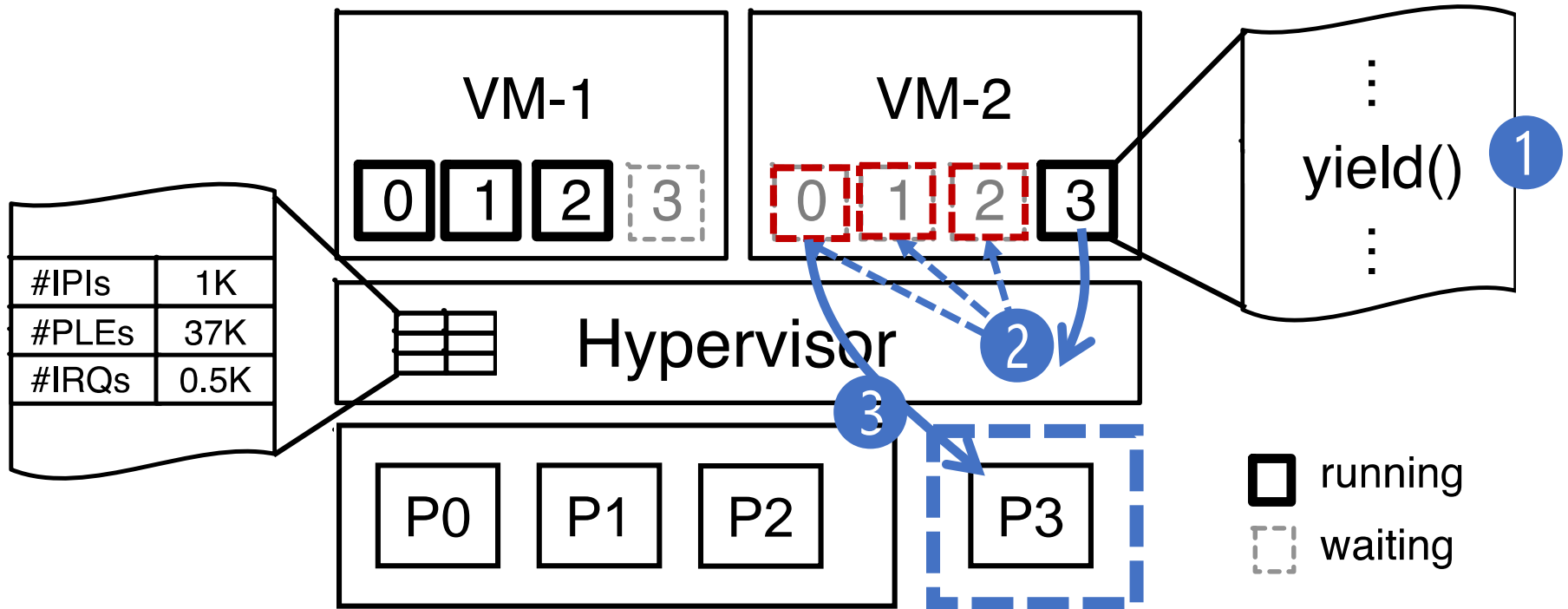
- 1 Yield occurring
- 2 Investigating the preempted vCPUs

Accelerating Critical TLB Synchronizations



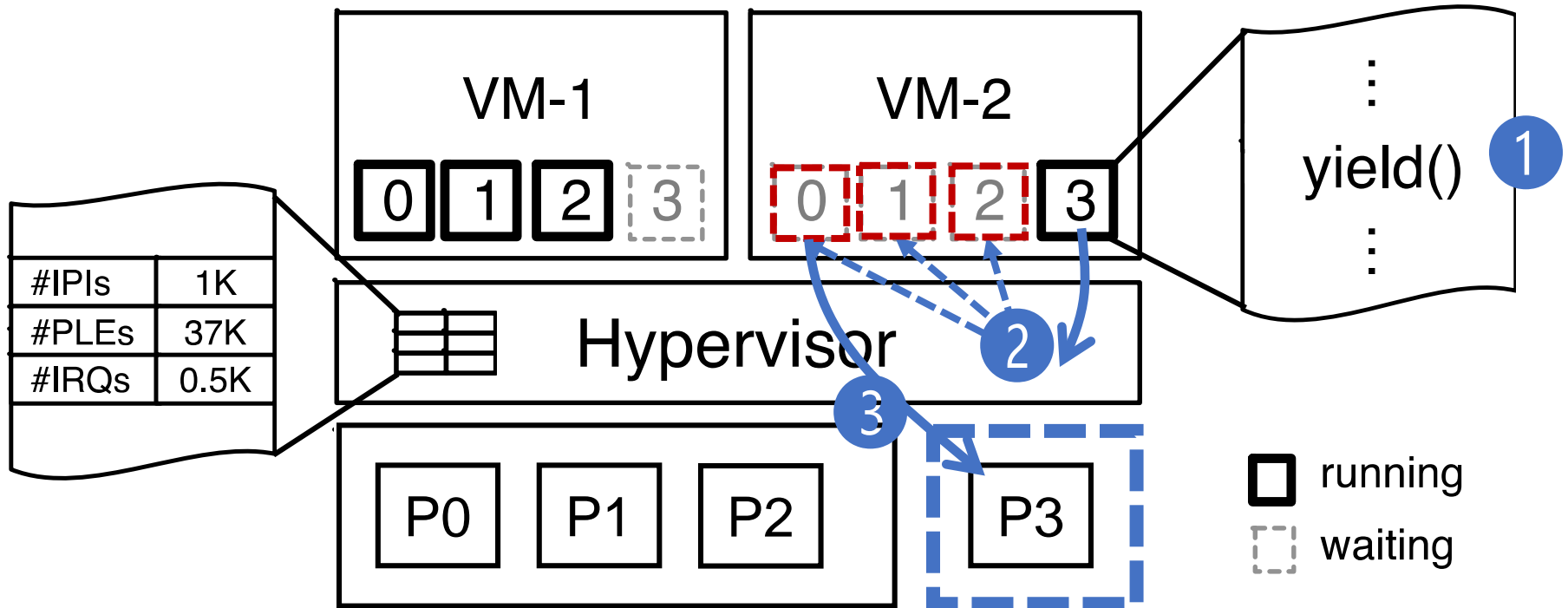
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical TLB Synchronizations



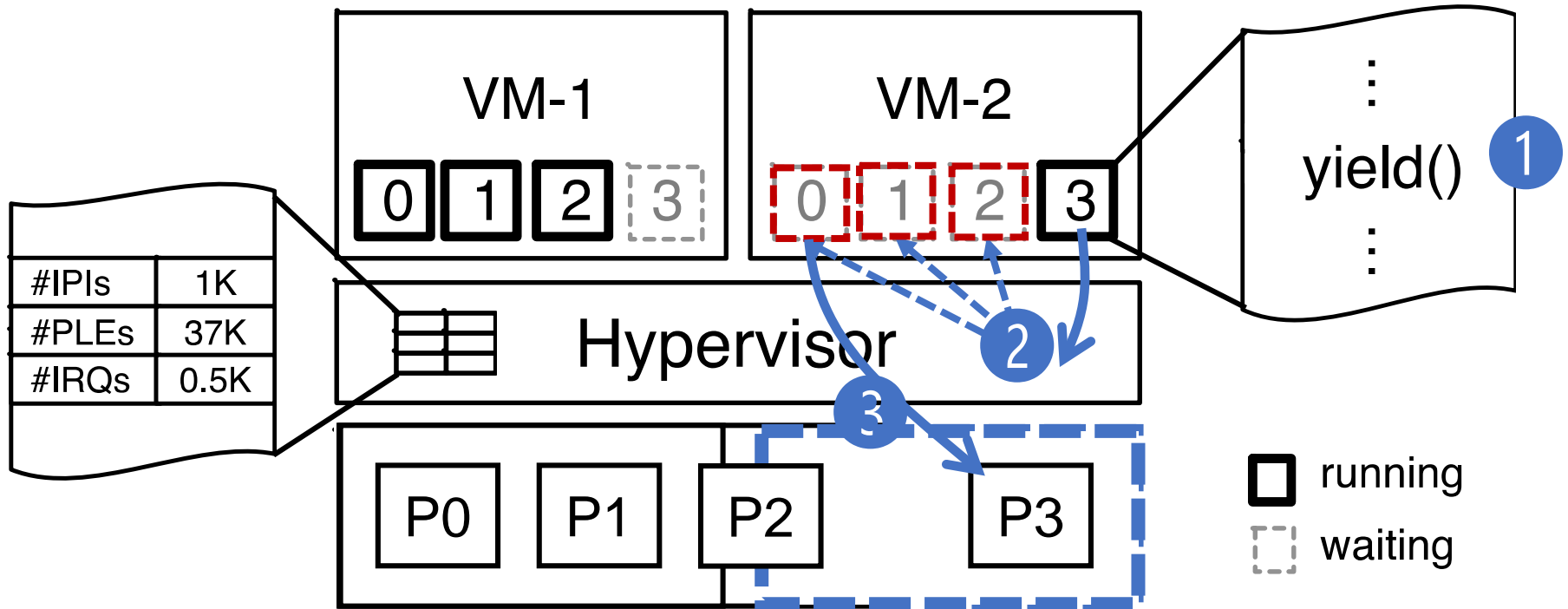
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool

Accelerating Critical TLB Synchronizations



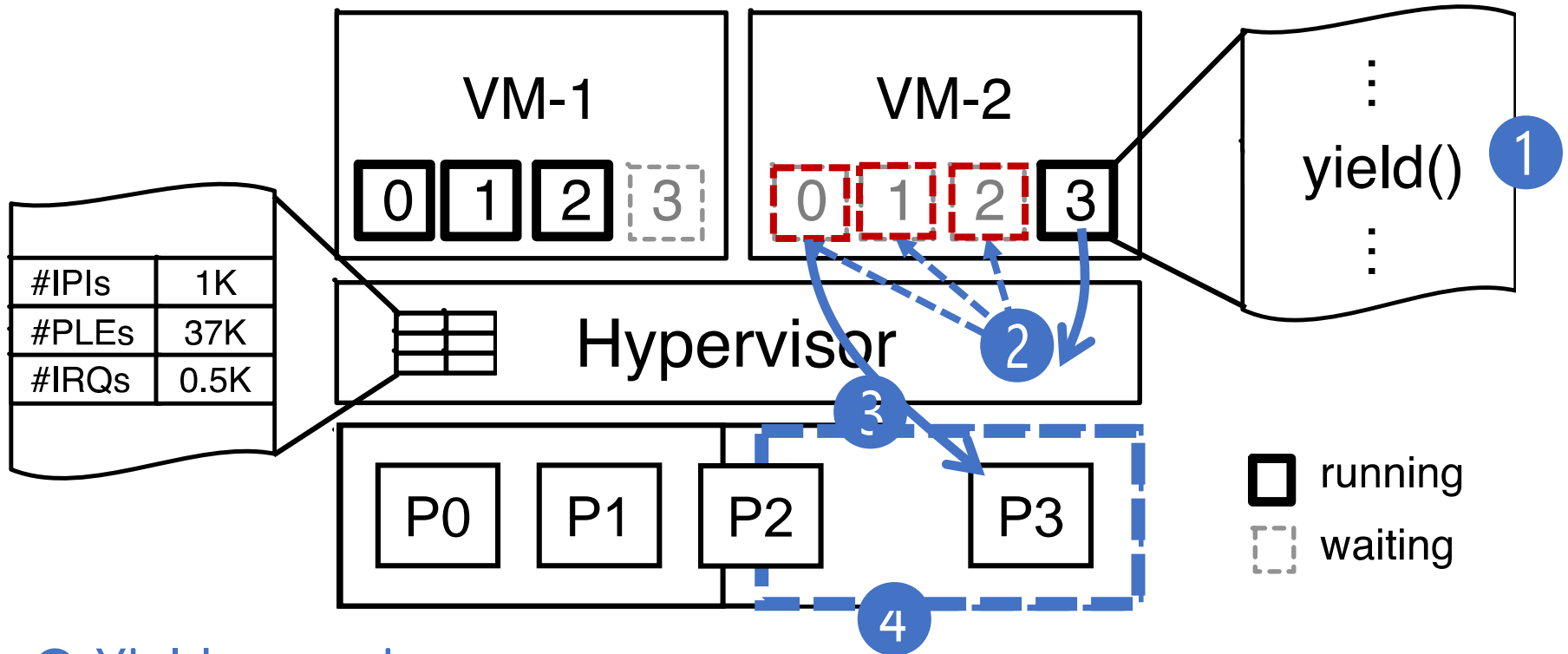
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



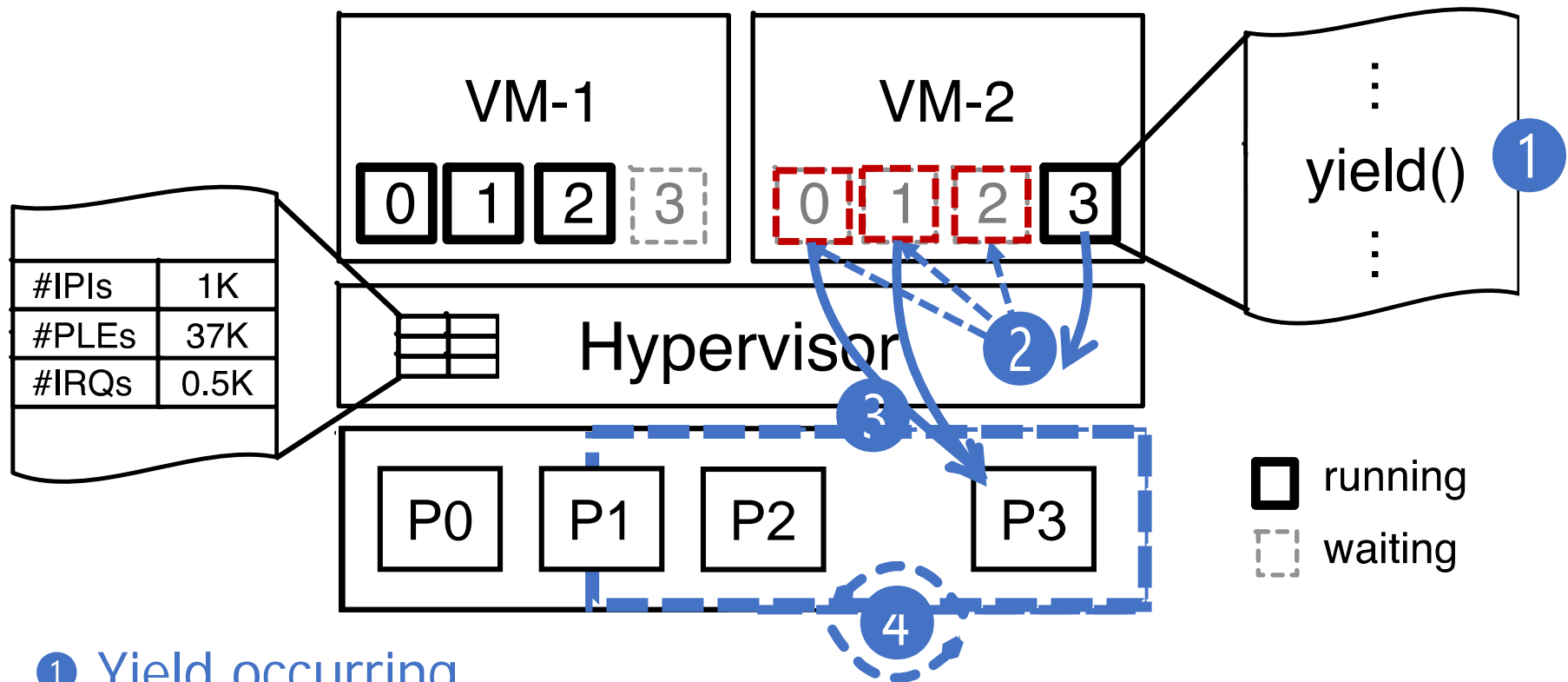
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



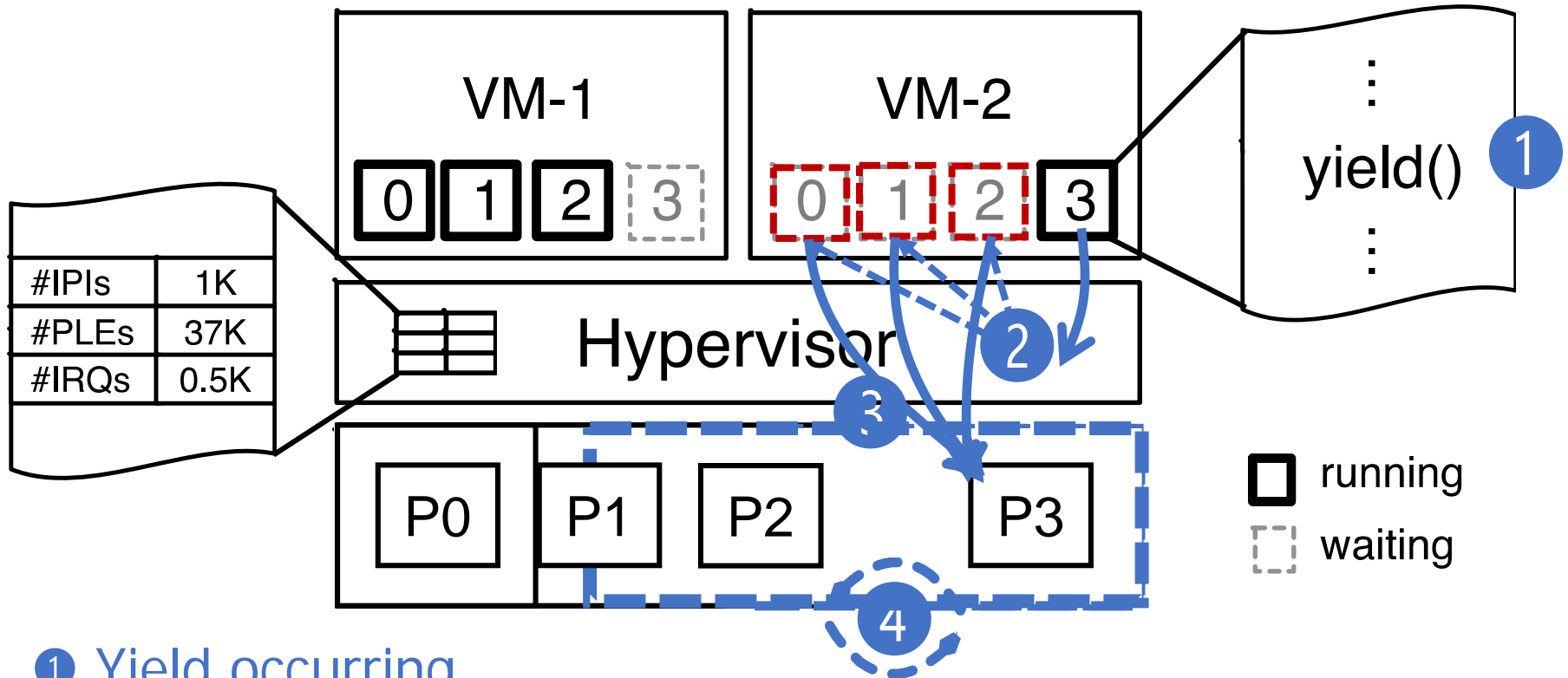
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



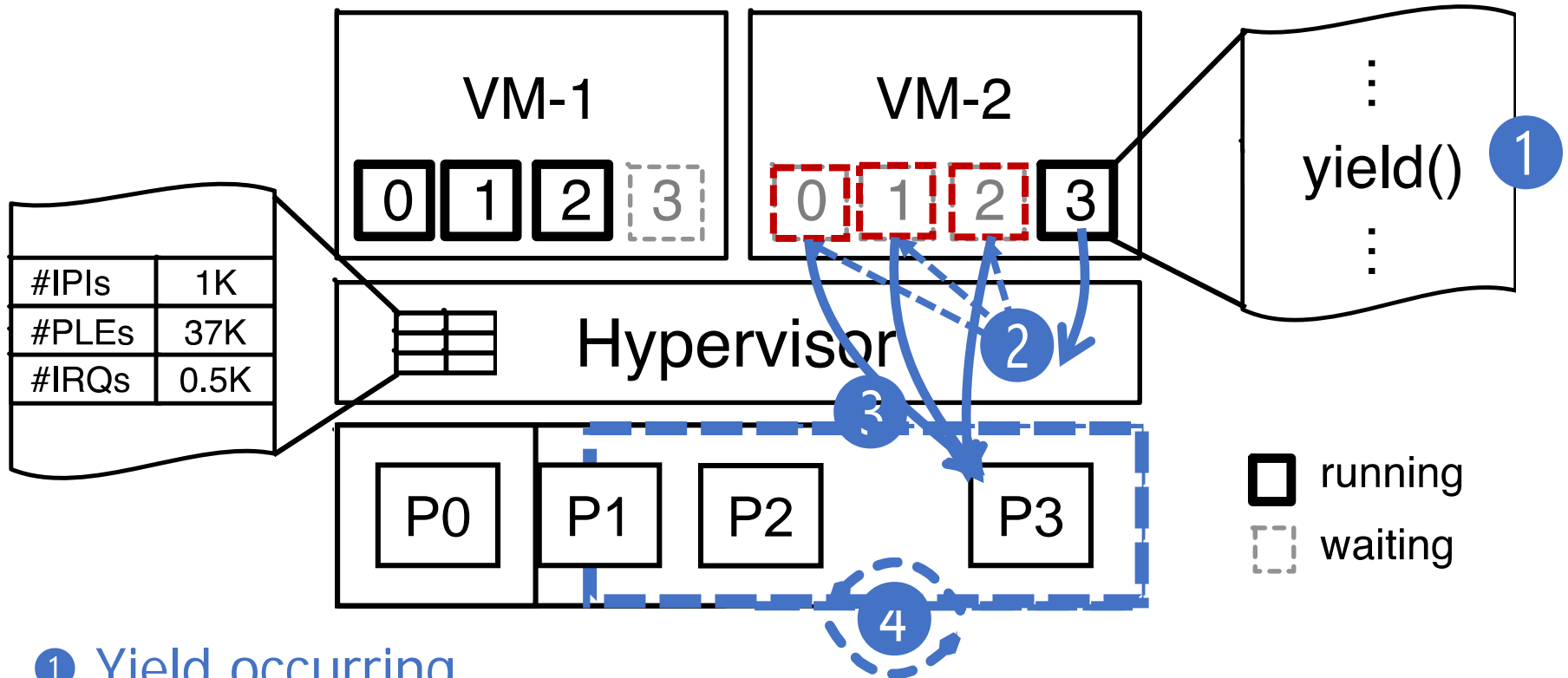
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



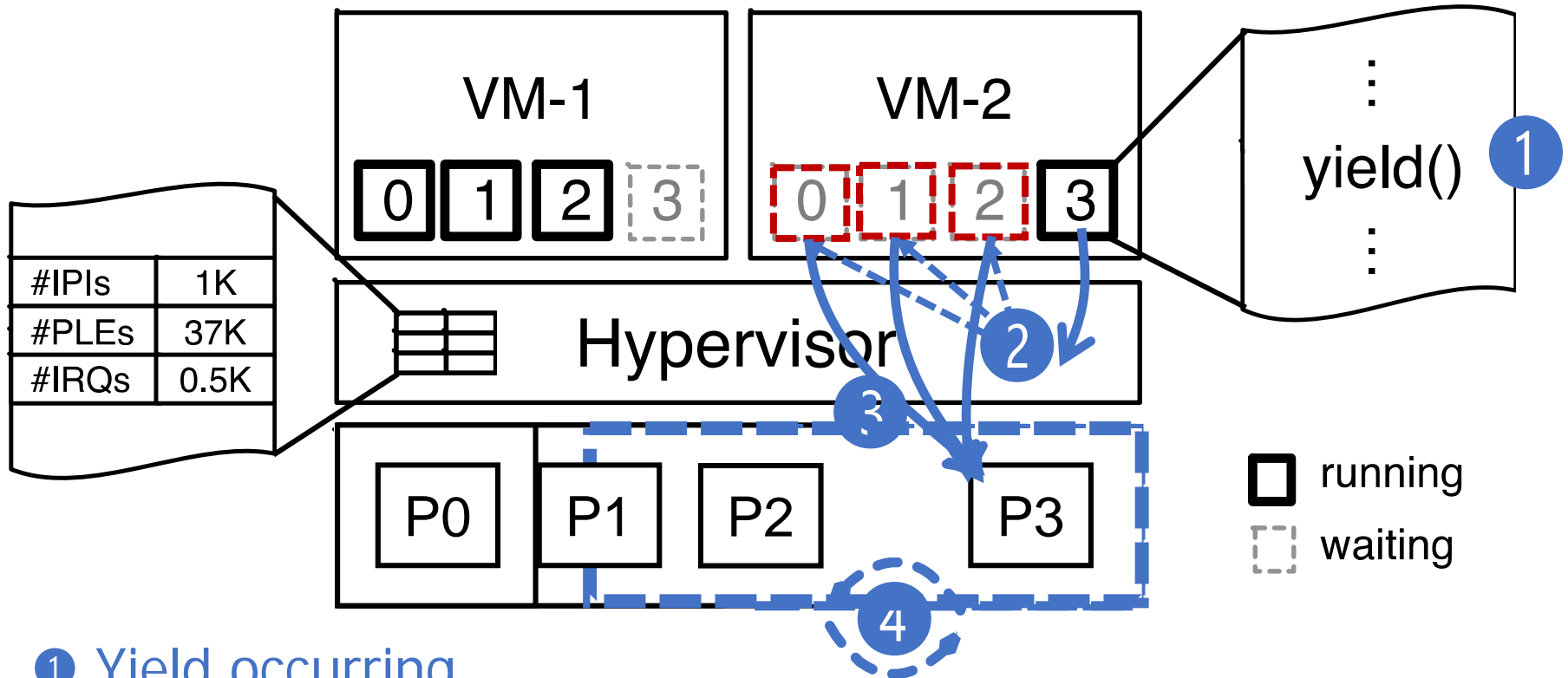
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



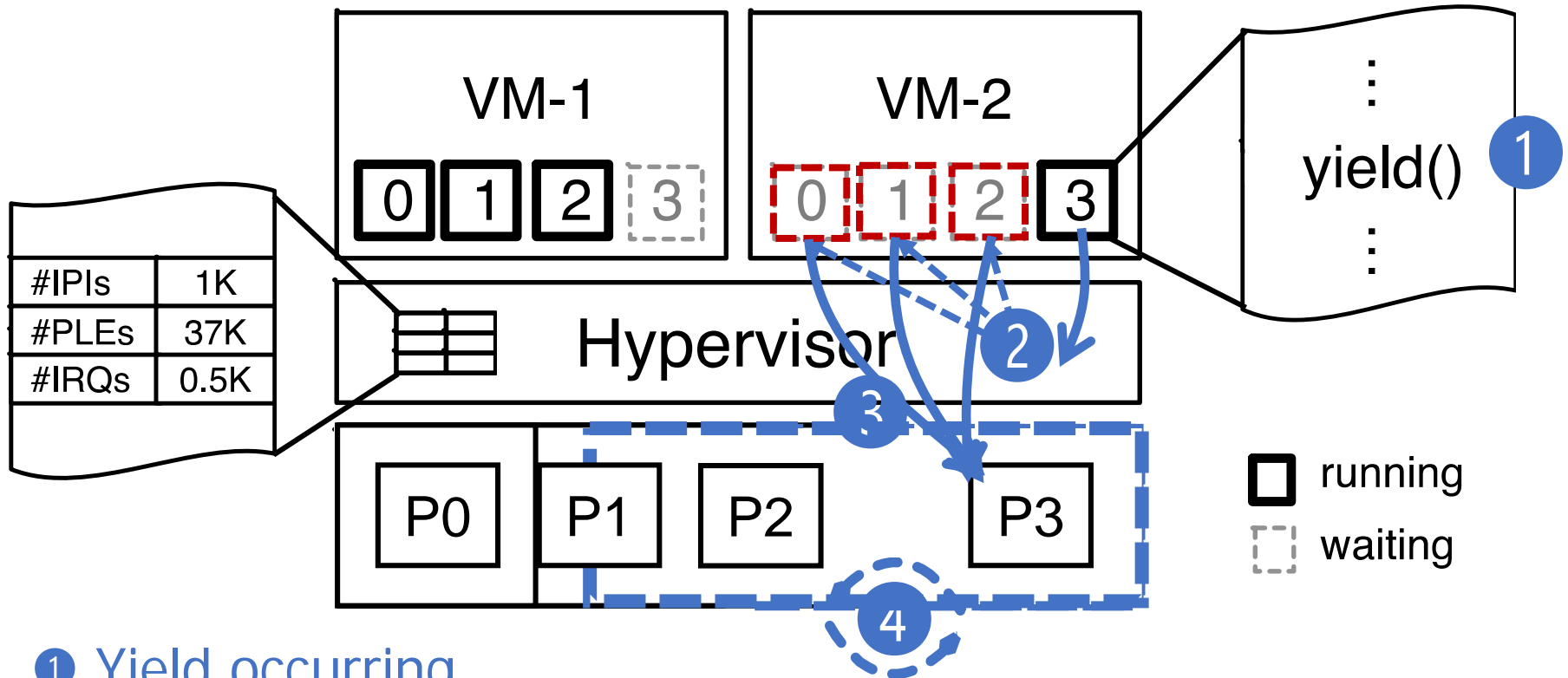
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



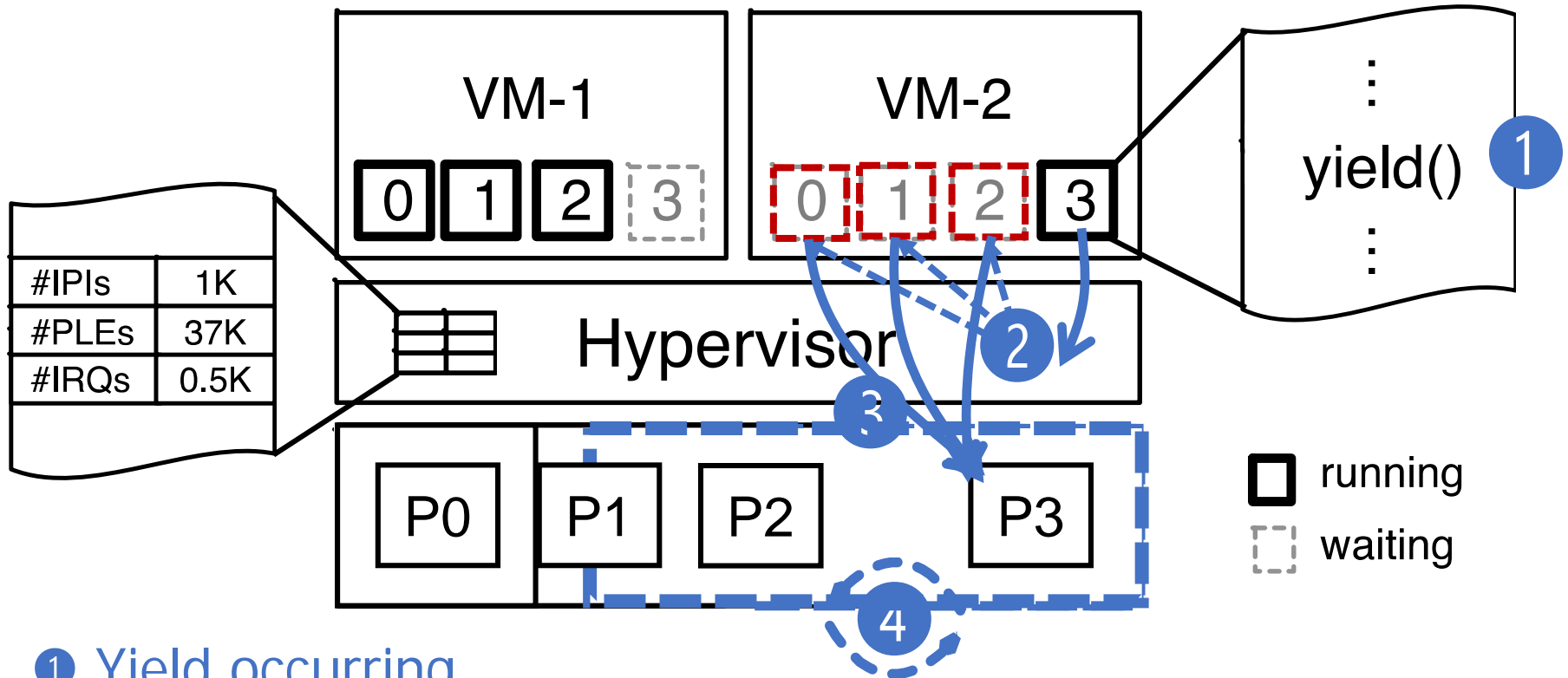
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



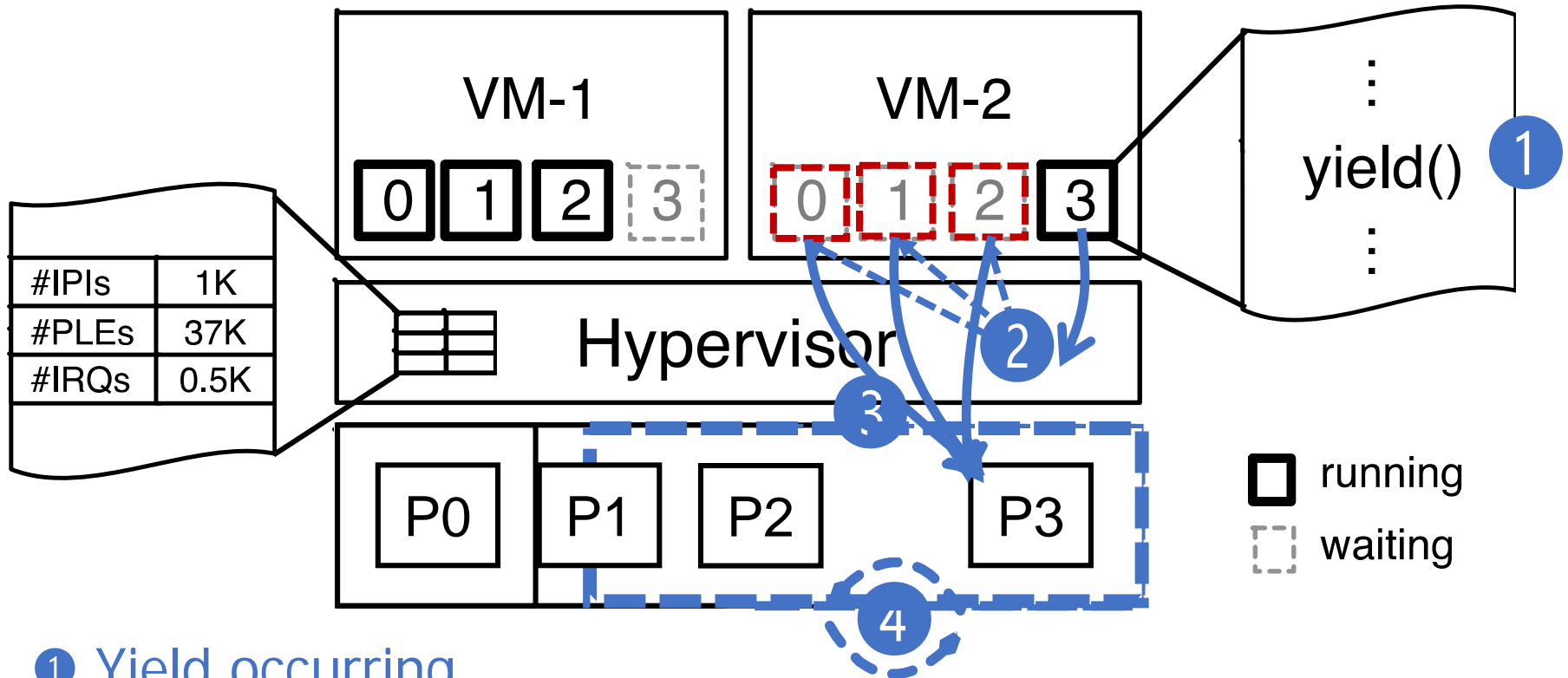
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



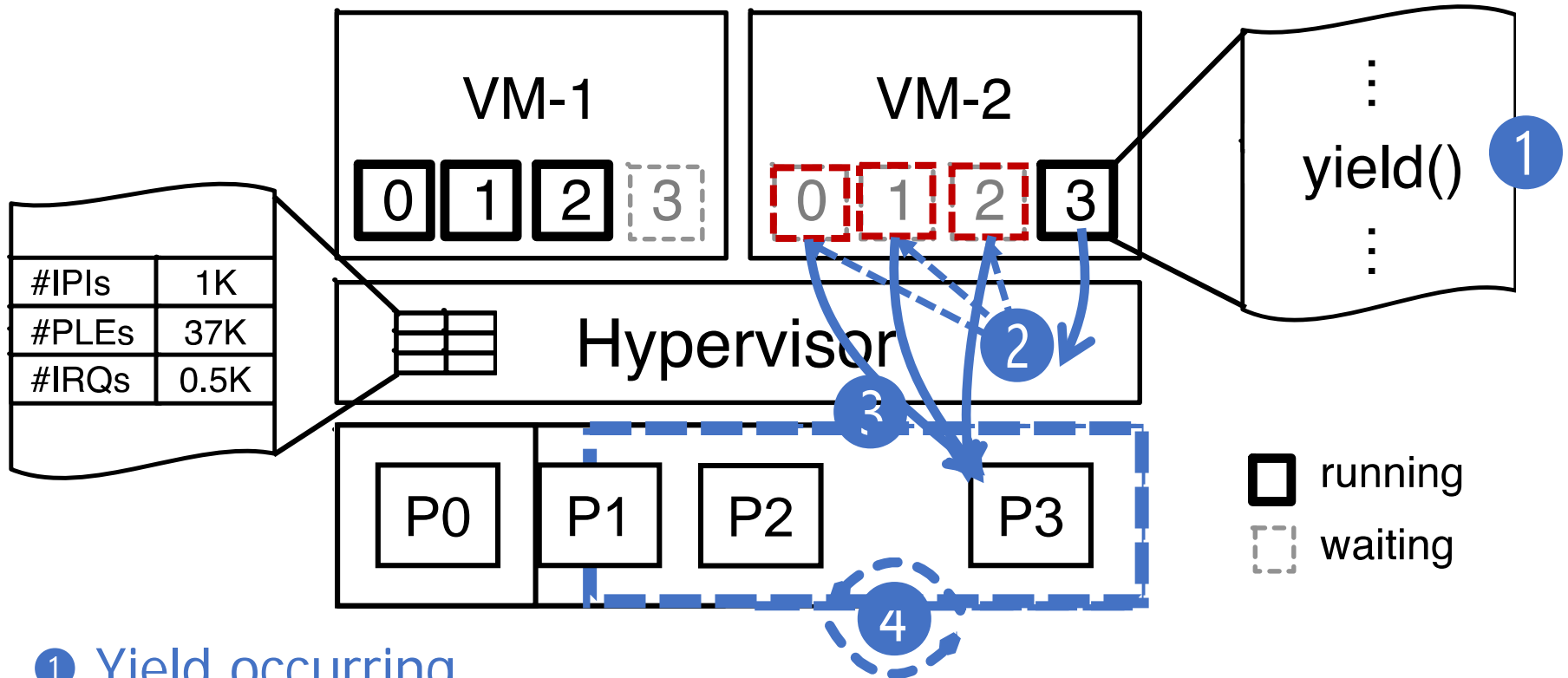
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



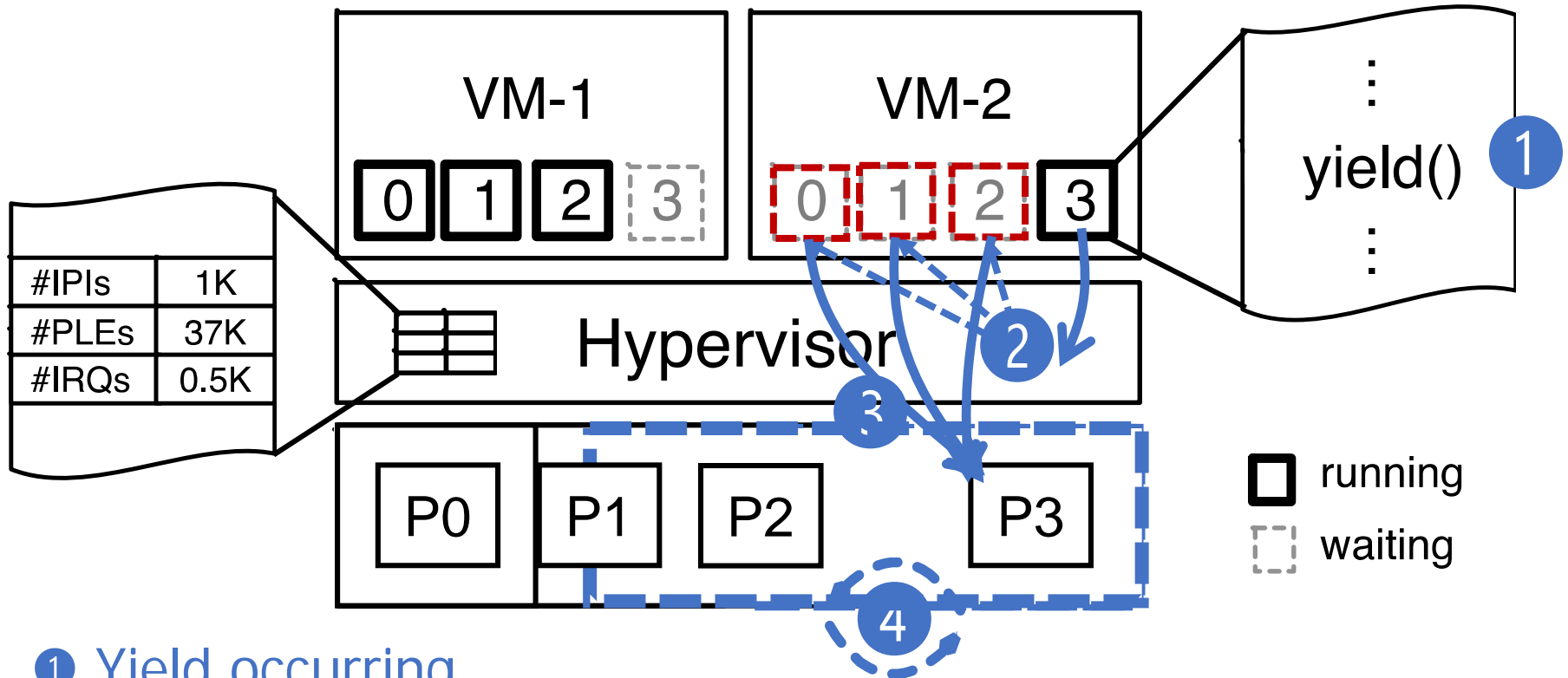
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



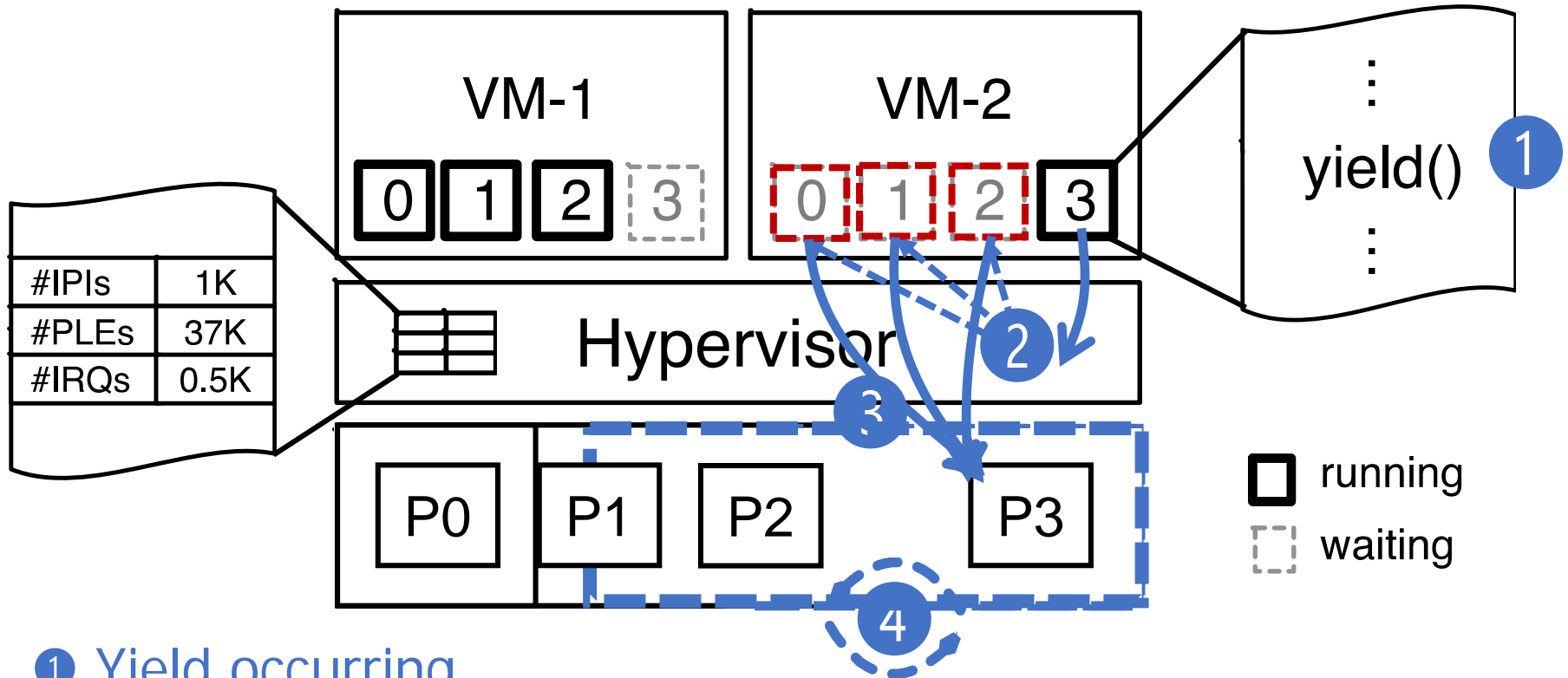
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



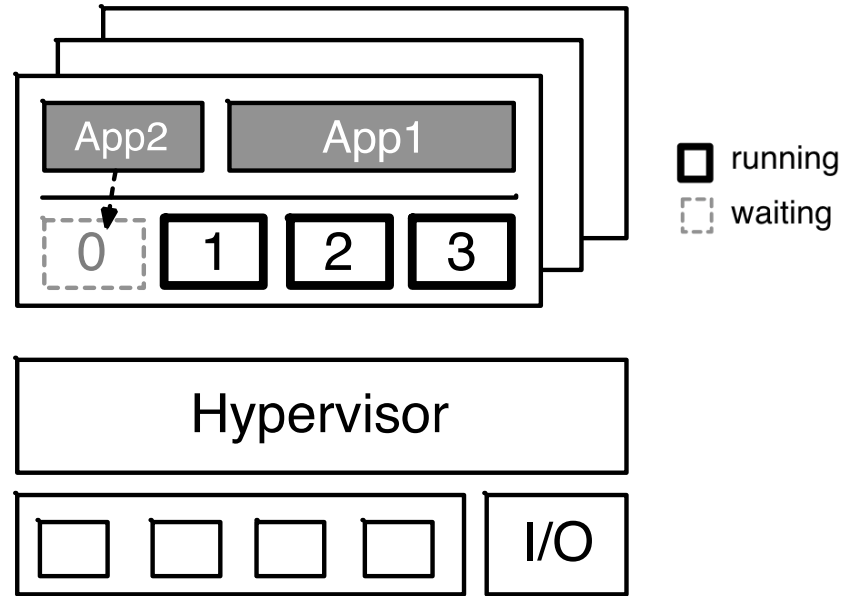
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Accelerating Critical TLB Synchronizations



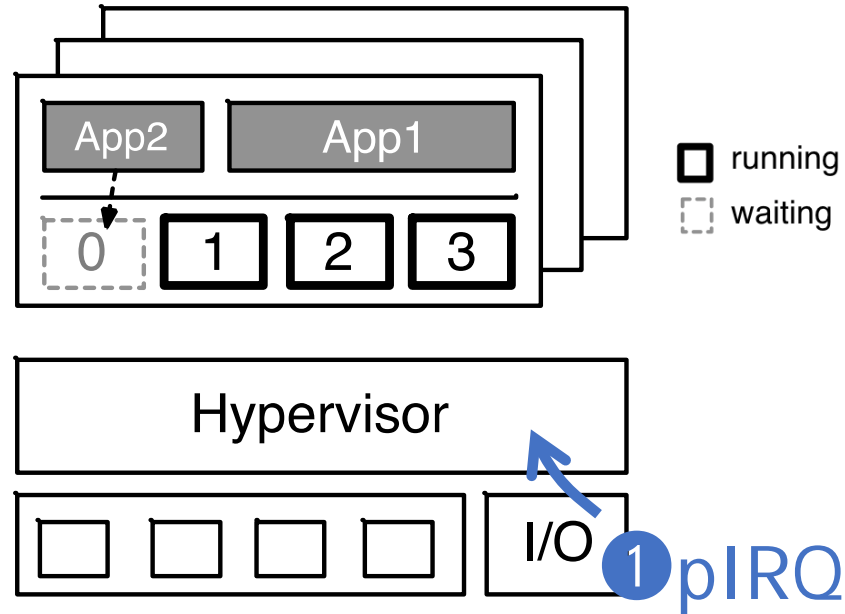
- 1 Yield occurring
- 2 Investigating the preempted vCPUs
- 3 Scheduling the selected vCPU on the micro-sliced pool
- 4 Dynamically adjusting micro-sliced cores based on profiling

Detecting Critical I/O Events



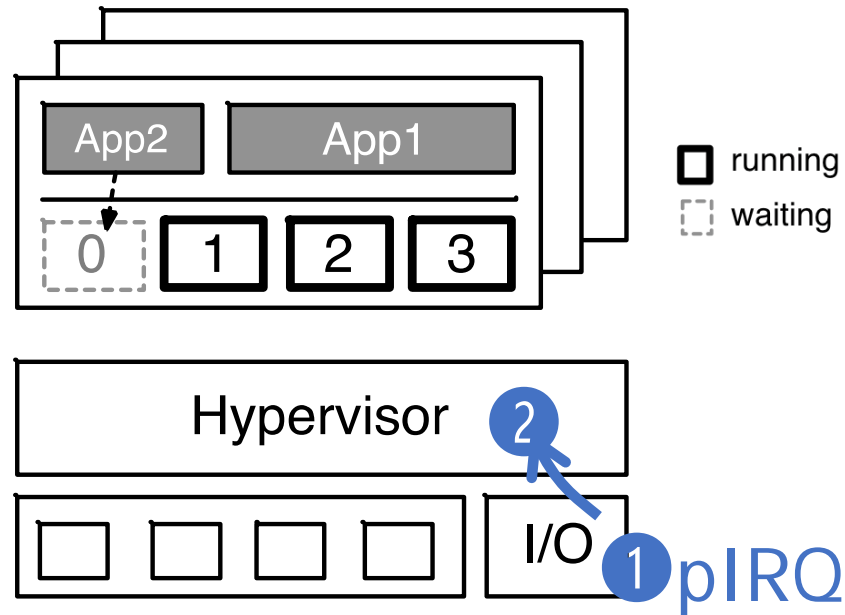
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



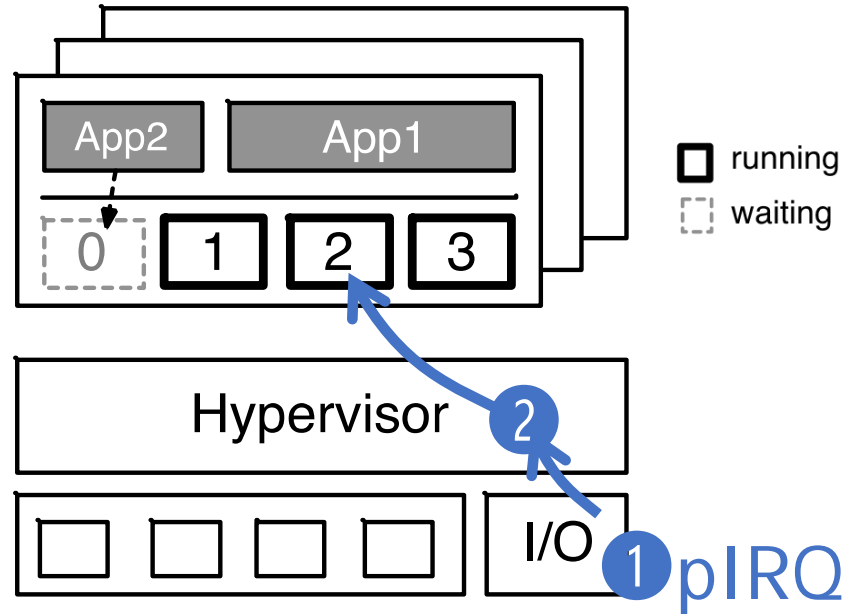
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



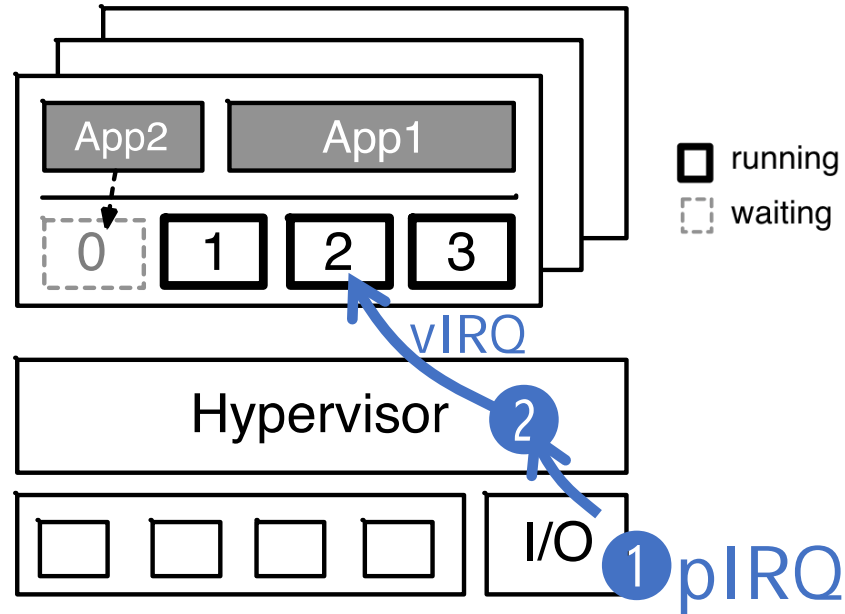
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



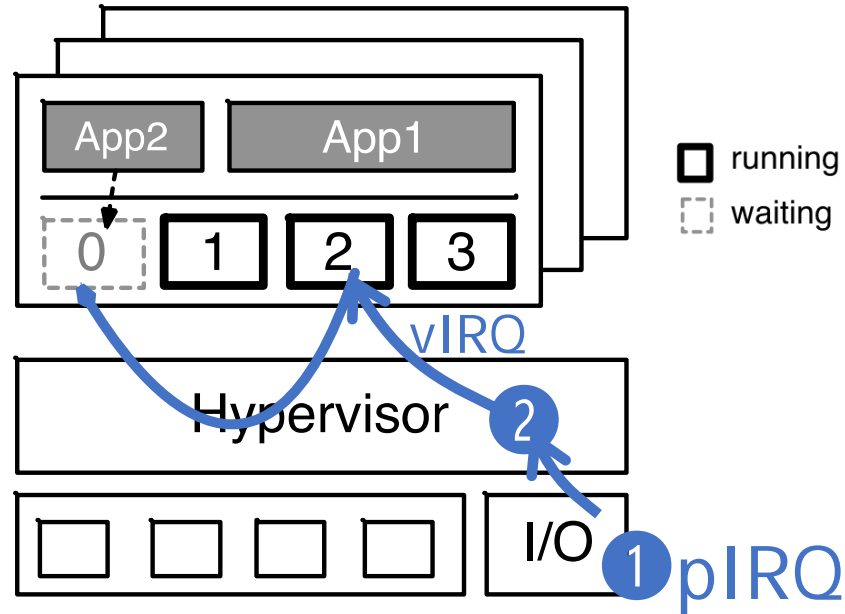
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



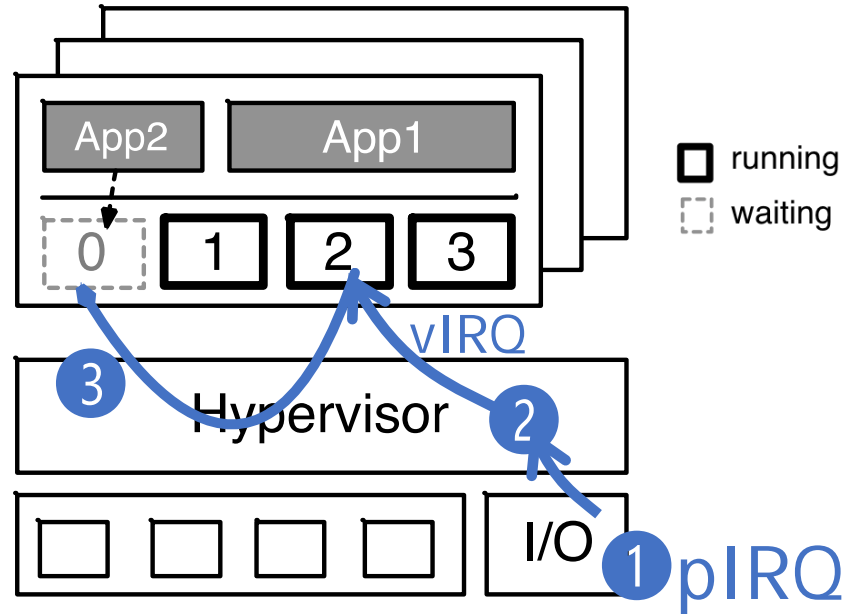
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



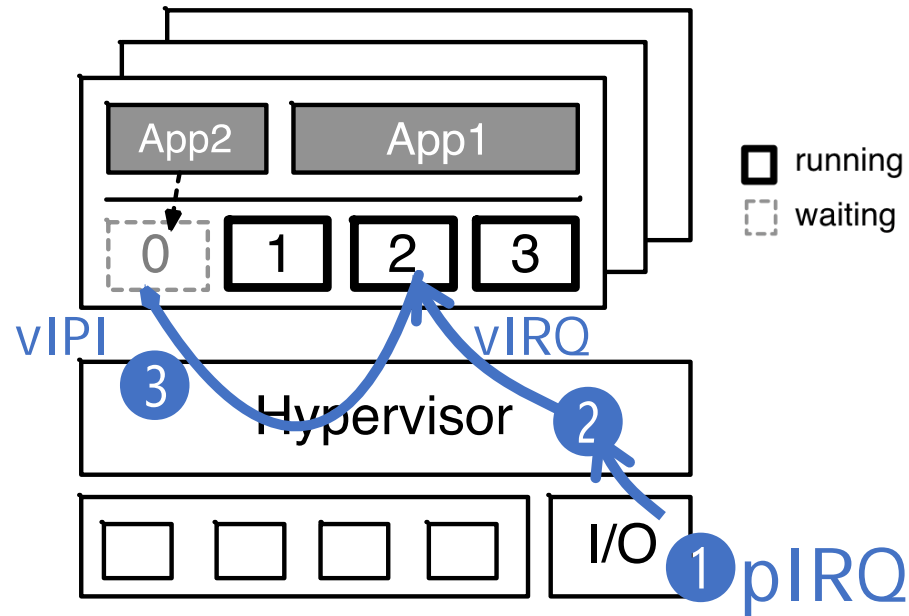
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



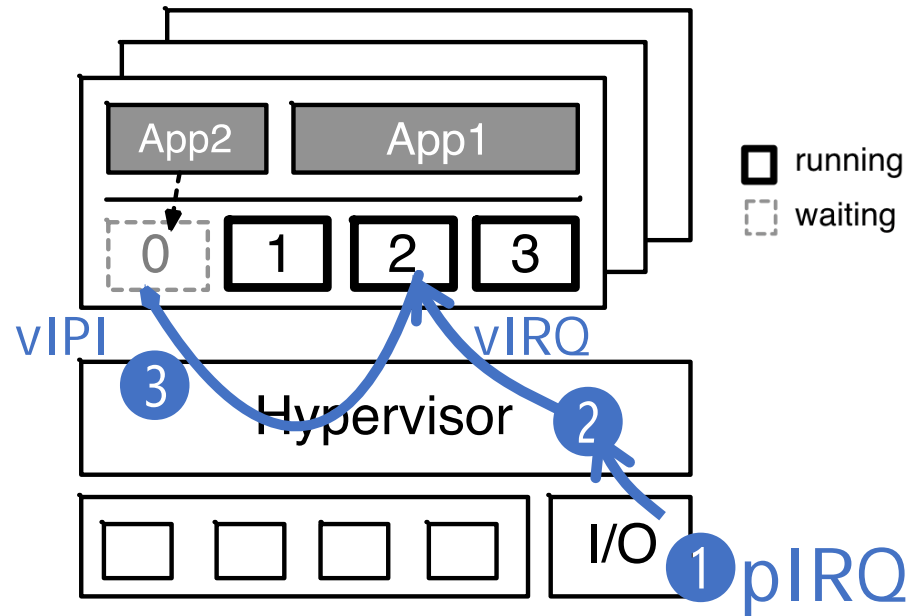
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Detecting Critical I/O Events



I/O handling consists of a chain of operations involving potentially multiple vCPUs

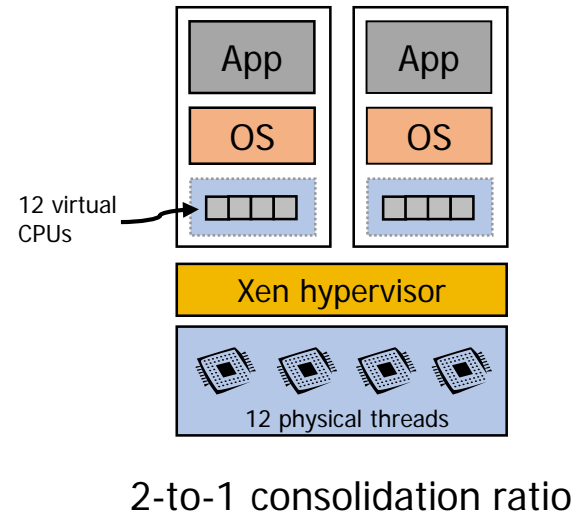
Detecting Critical I/O Events



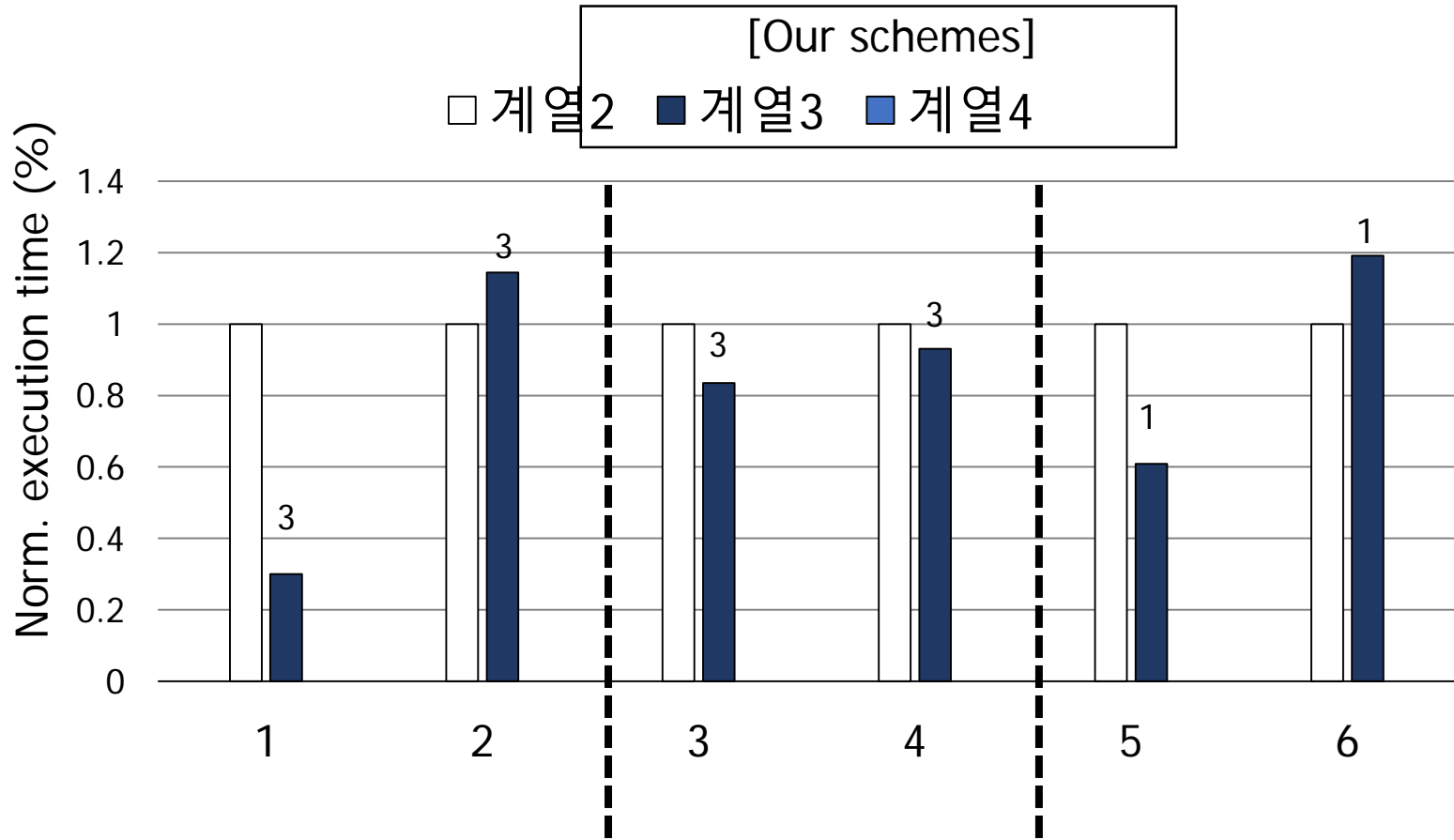
I/O handling consists of a chain of operations involving potentially multiple vCPUs

Experimental Environments

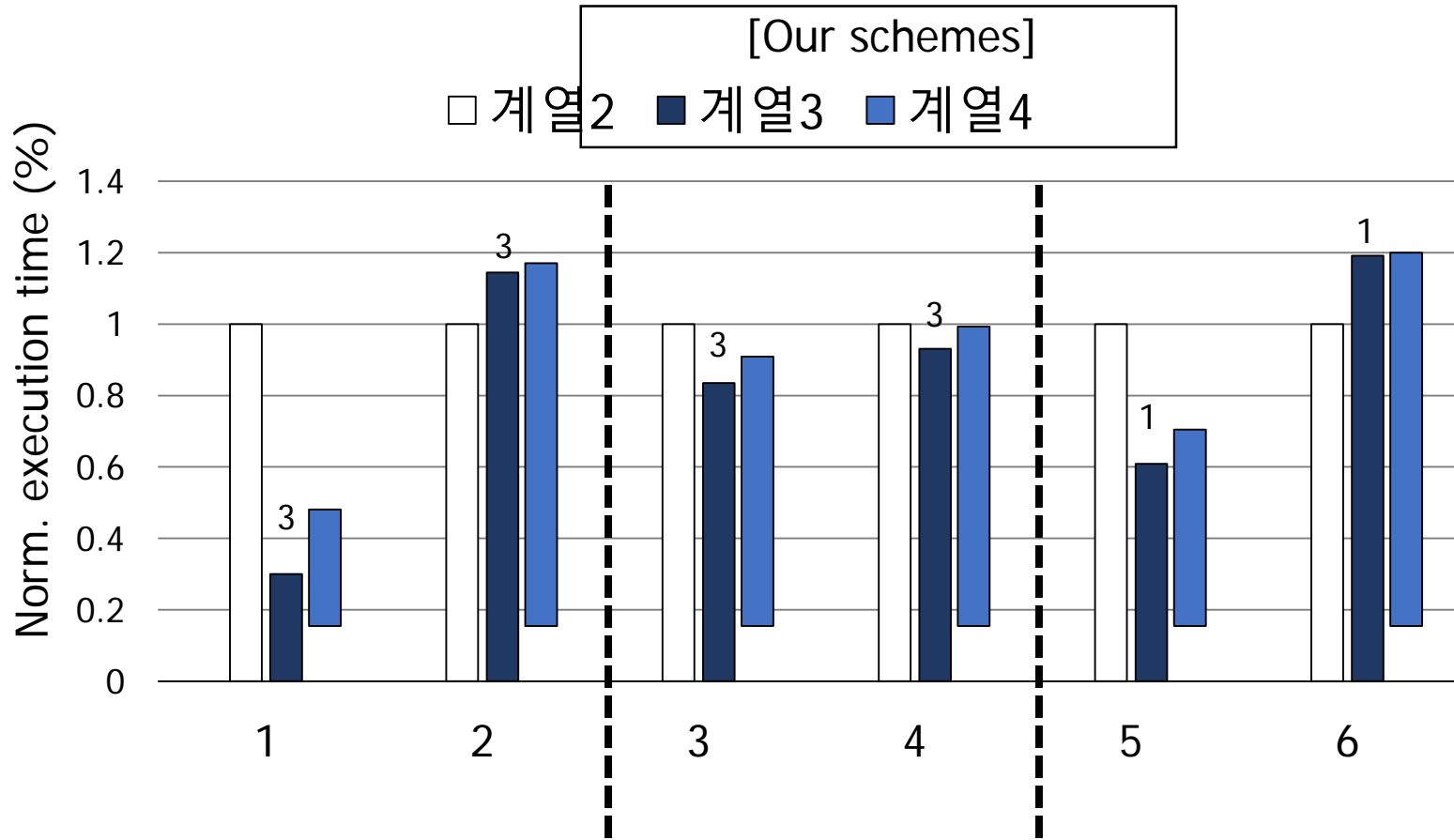
- Testbed
 - 12 HW threads (Intel Xeon)
 - 2 VMs with 12 vCPUs for each
 - Xen hypervisor 4.7
- Benchmarking workloads
 - `dedup` and `vips` from PARSEC
 - `exim` and `gmake` from MOSBENCH
- Pool configuration
 - Normal: 30ms (Xen default)
 - Micro-sliced: 0.1ms



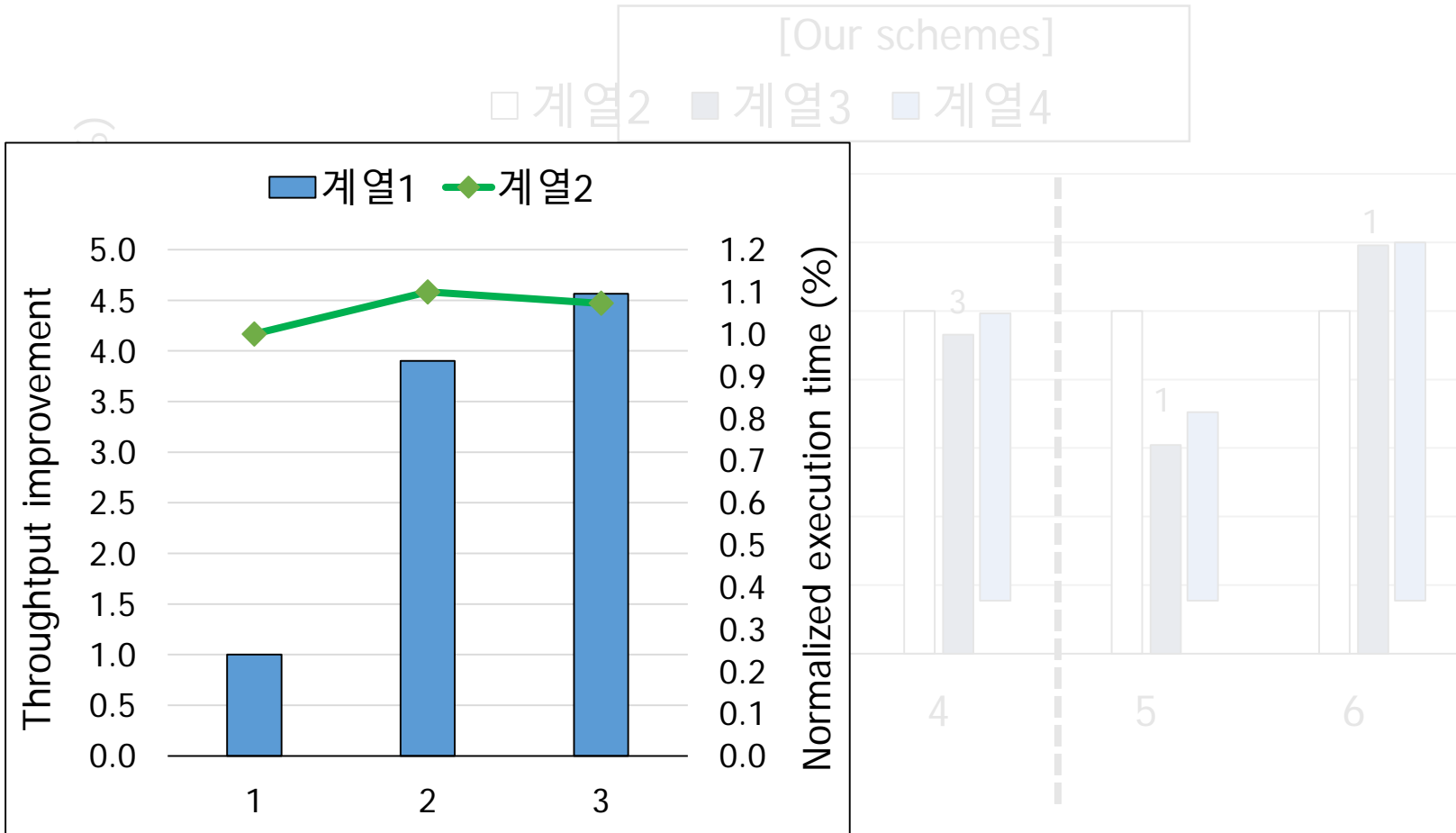
Performance of Micro-sliced Cores



Performance of Micro-sliced Cores



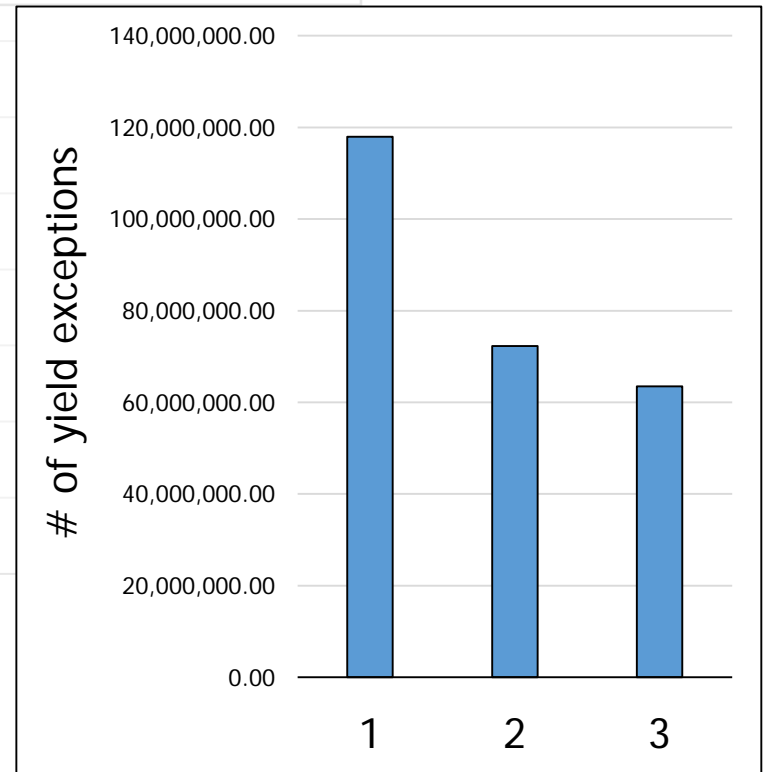
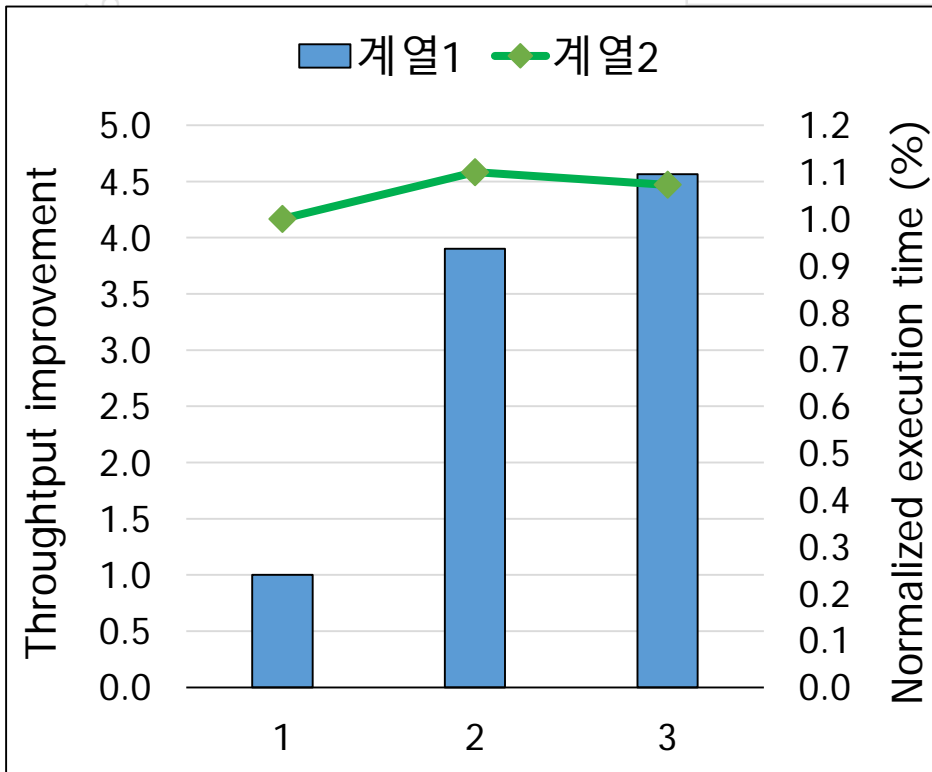
Performance of Micro-sliced Cores



Performance of Micro-sliced Cores

[Our schemes]

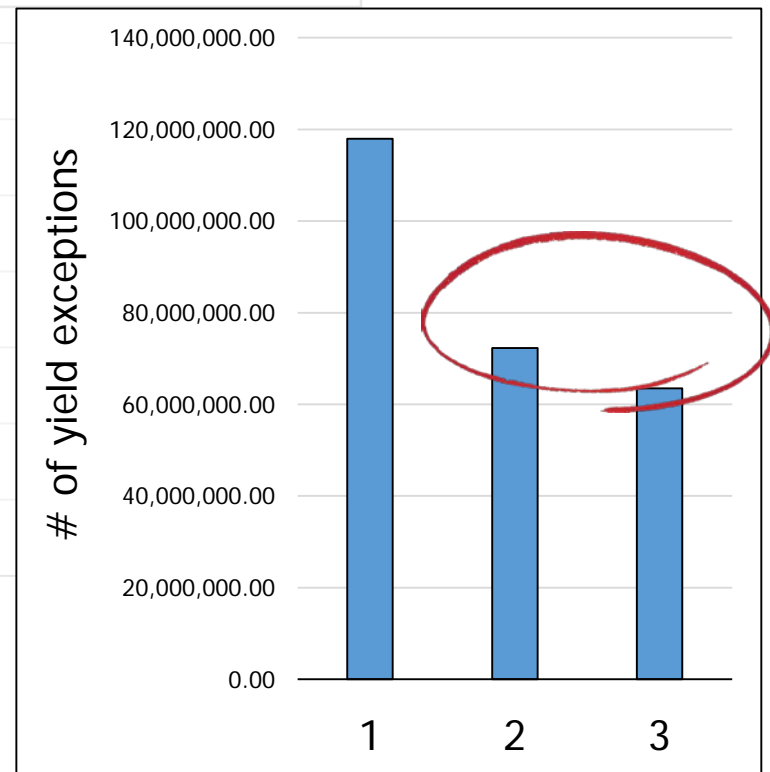
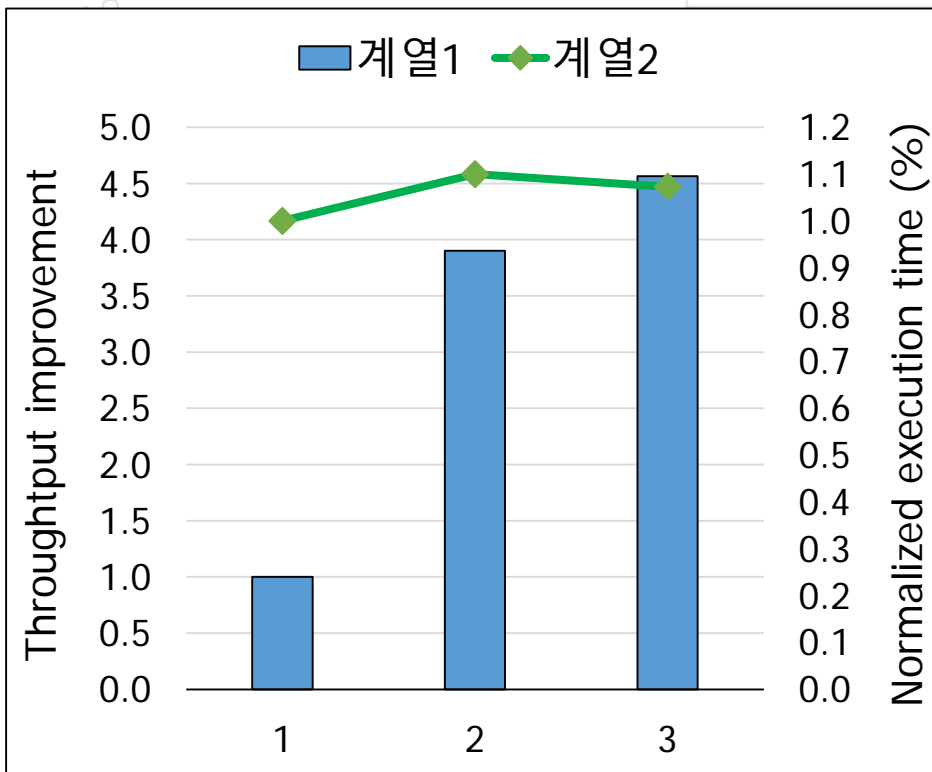
□ 계열2 ■ 계열3 ■ 계열4



Performance of Micro-sliced Cores

[Our schemes]

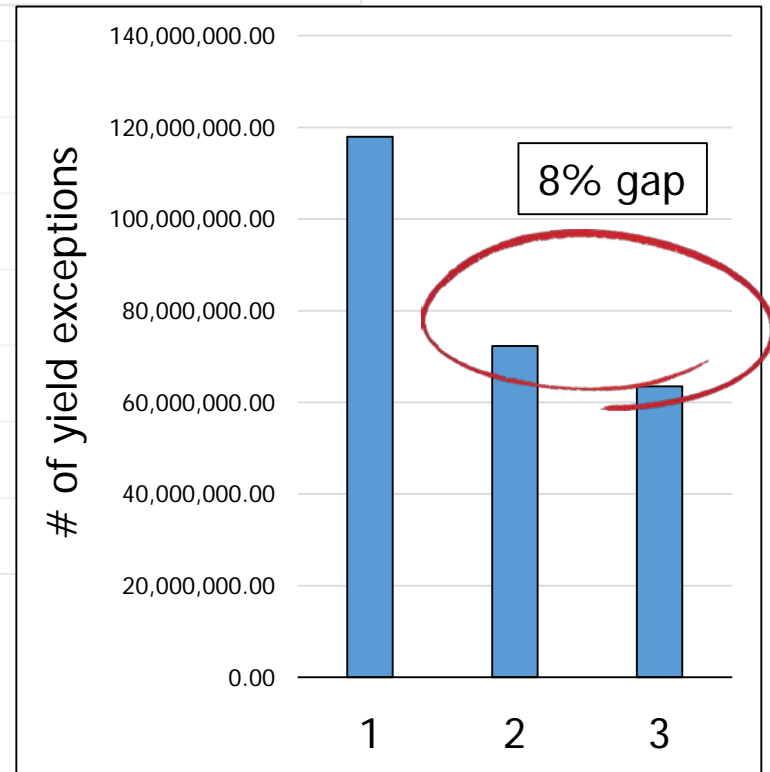
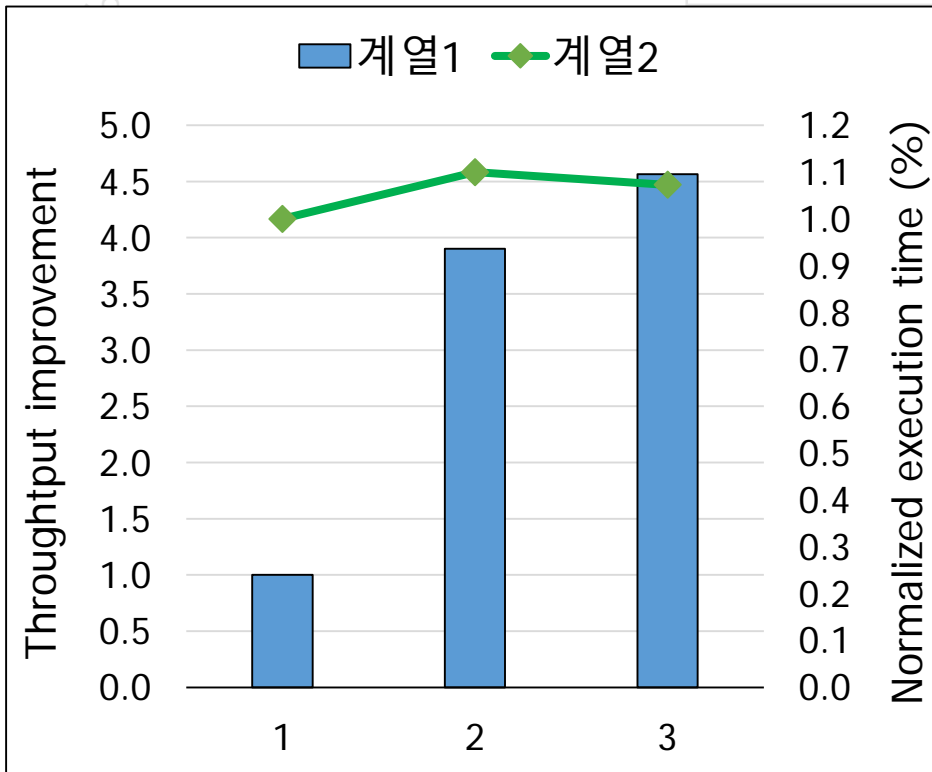
□ 계열2 ■ 계열3 ■ 계열4



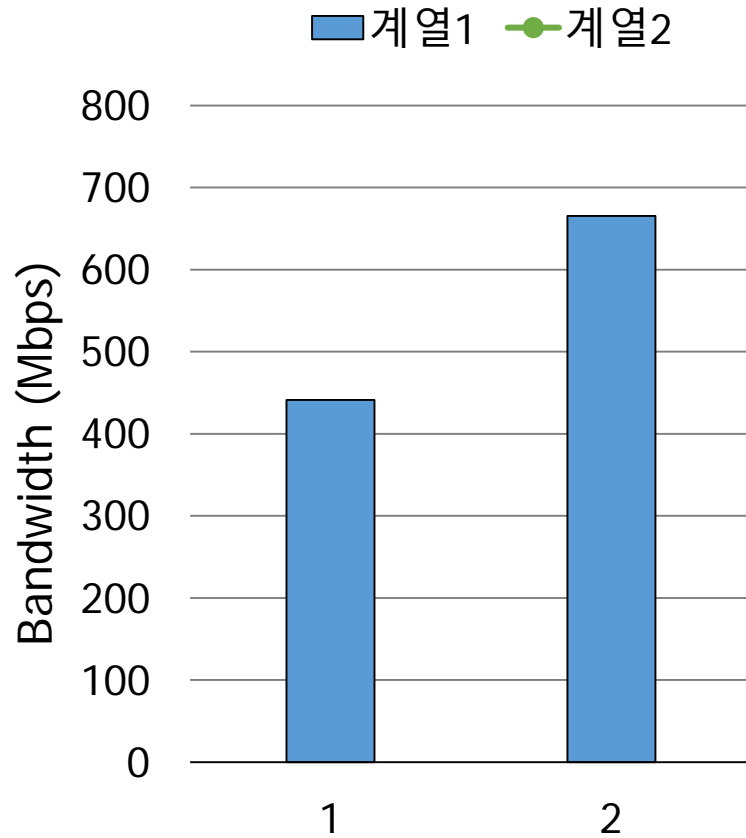
Performance of Micro-sliced Cores

[Our schemes]

□ 계열2 ■ 계열3 ■ 계열4

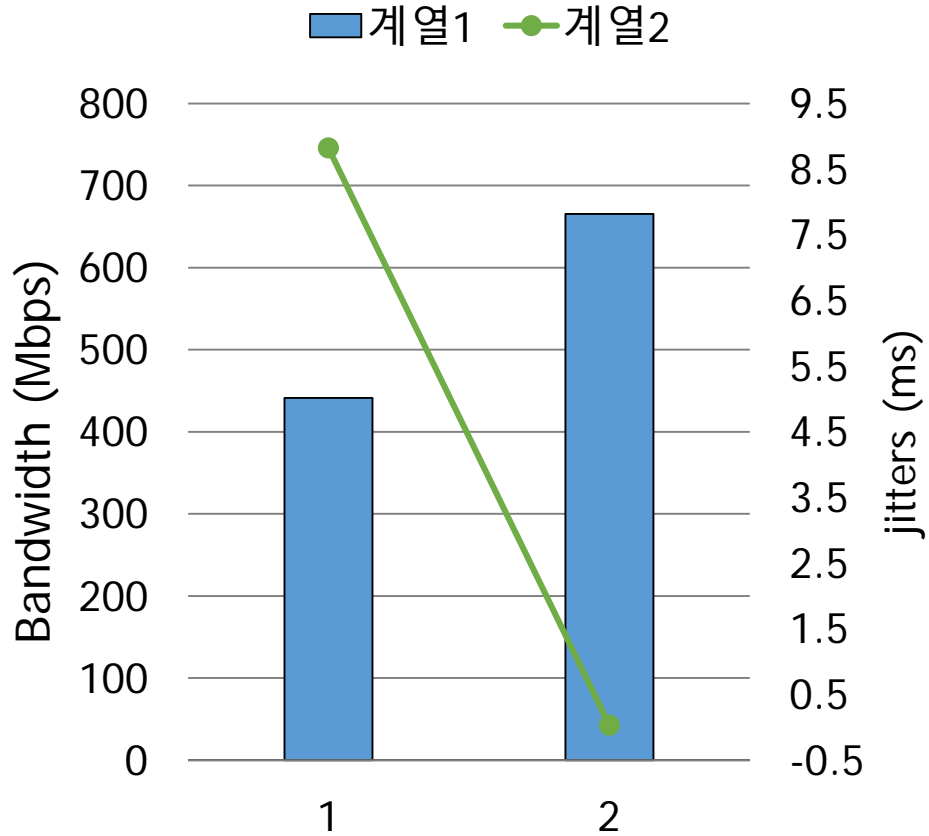


I/O Performance



	Workloads
VM-1	iPerf lookbusy
VM-2	lookbusy

I/O Performance



	Workloads
VM-1	iPerf lookbusy
VM-2	lookbusy

Conclusions

- We introduced a new approach to mitigate the virtual time discontinuity problem
- Three distinct contributions
 - Precise detection of urgent tasks
 - Guest OS transparency
 - Dynamic adjustment of the micro-sliced cores
- Overhead is very low for applications which do not frequently use OS services

Thank You!

jsahn@ajou.ac.kr

Accelerating Critical OS Services in Virtualized
Systems with Flexible Micro-sliced Cores

Jeongseob Ahn^{*}, Chang Hyun Park[‡], Taekyung Heo[‡], Jaehyuk Huh[‡]

*



AJOU UNIVERSITY

‡

KAIST

