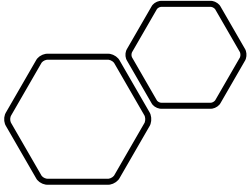


73,000 Pods a Day Misadventures in Multi-Tenant

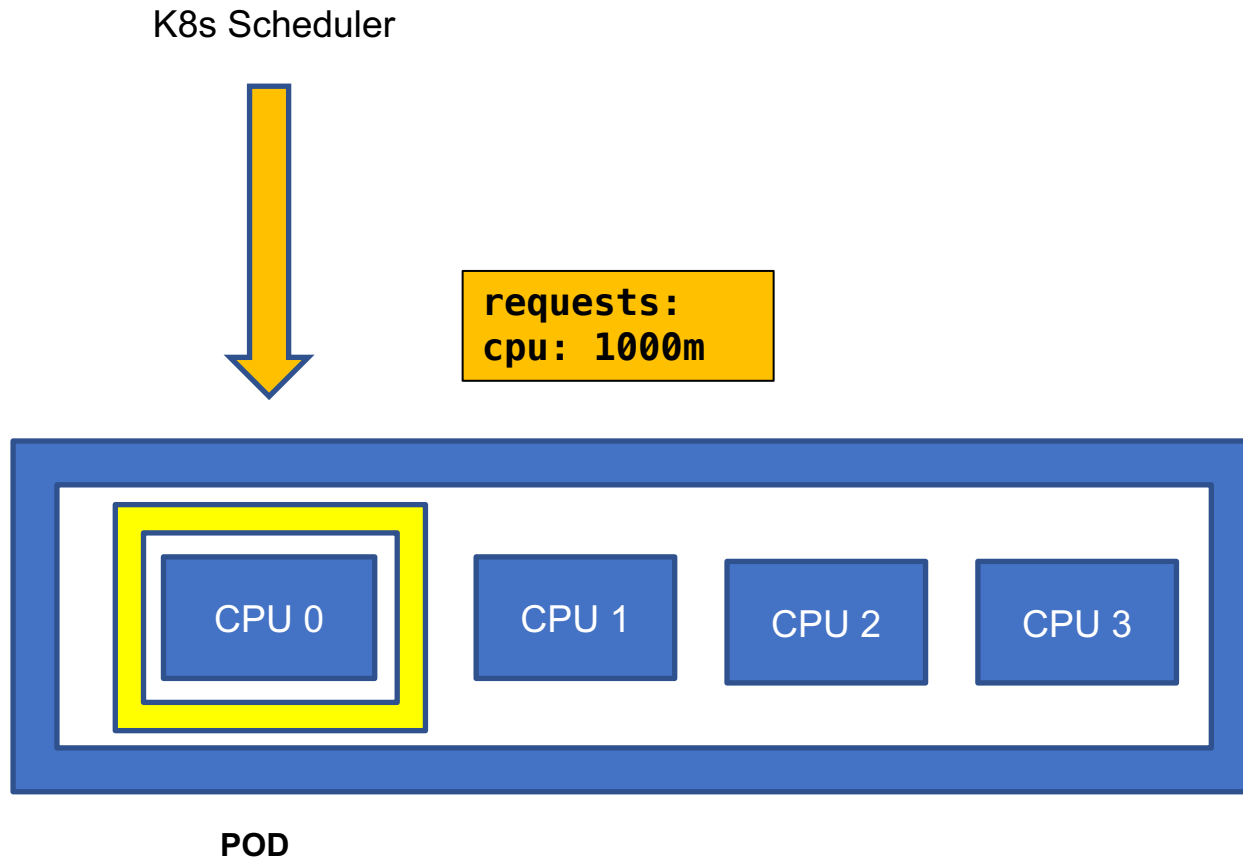
Wil Reed, Acquia

Shane Corbett, Amazon Web Services



Kubernetes vs. Linux

Cores

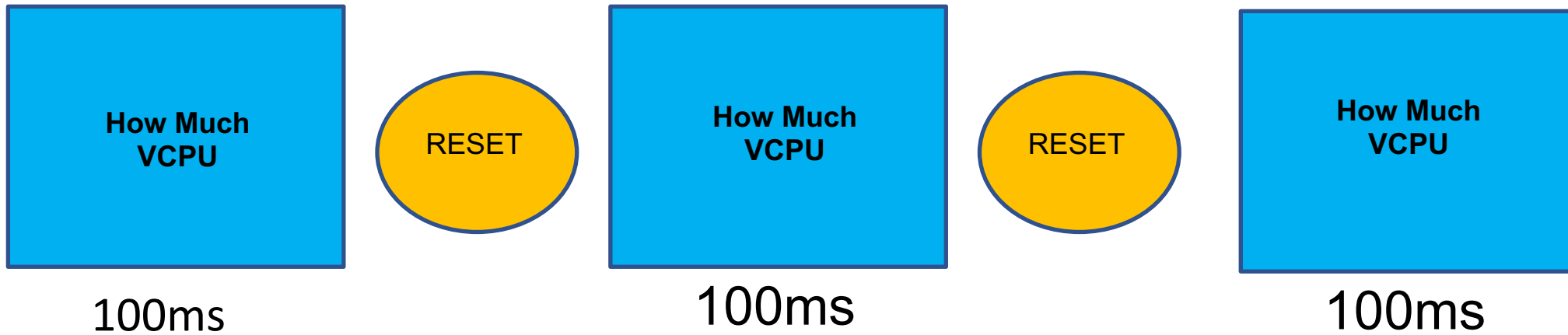


Capacity:	
cpu:	4
Pods:	110
Allocatable:	
cpu:	3000m
Pods:	110



**Thinking in Time
Not Cores**

Limits



Period – 100ms

Conversion

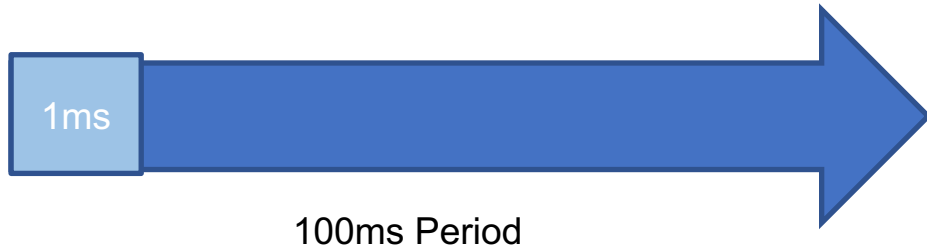
1 Core = 1000m

Resources:

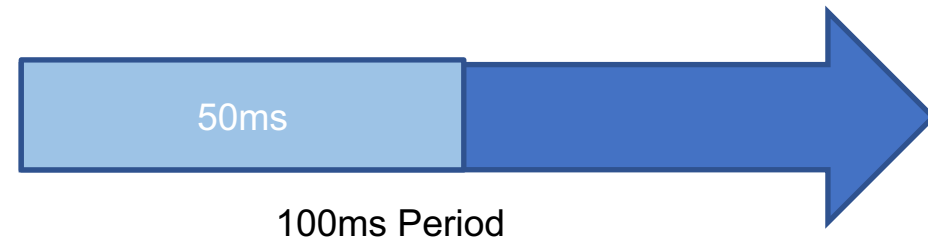
limits:

cpu: 10m

$$10\text{m}/1000\text{m} = 1\%$$



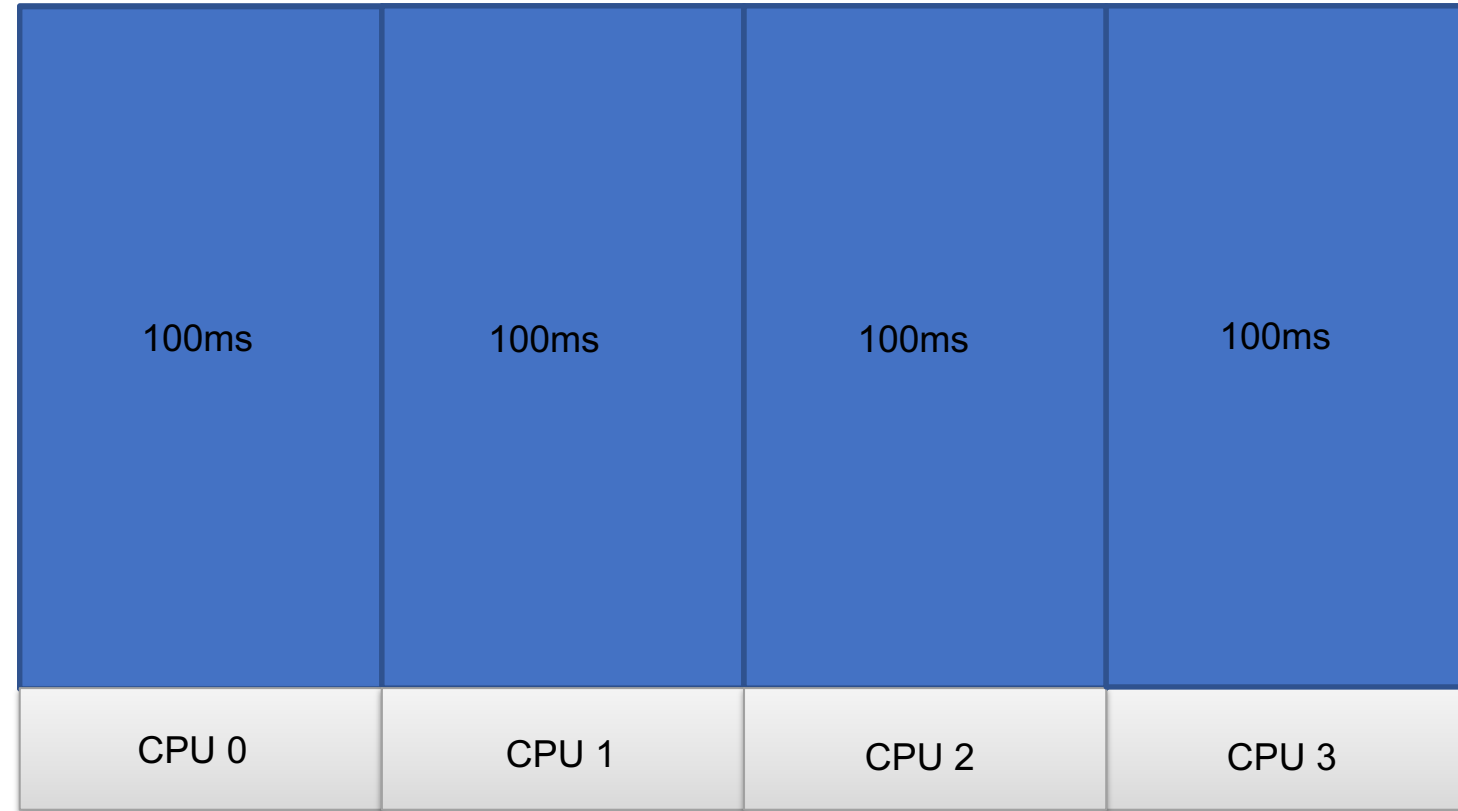
$$500\text{m}/1000\text{m} = 50\%$$



Time

Four Thread app

4 Core Node

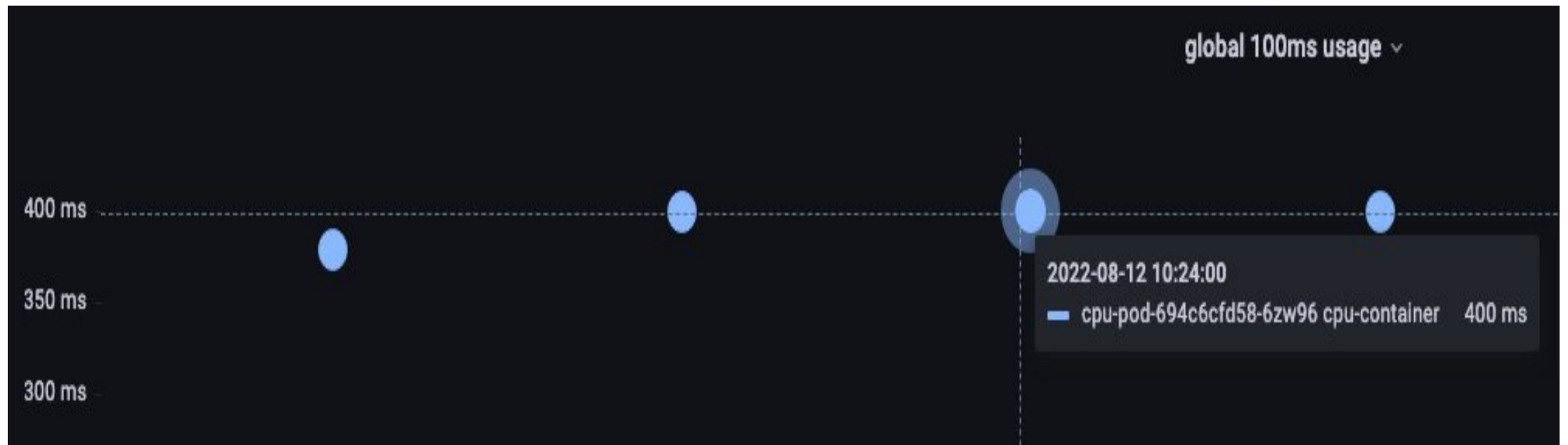


Real Time
100ms

Quota Time
400ms = 4000m

Thinking in Time

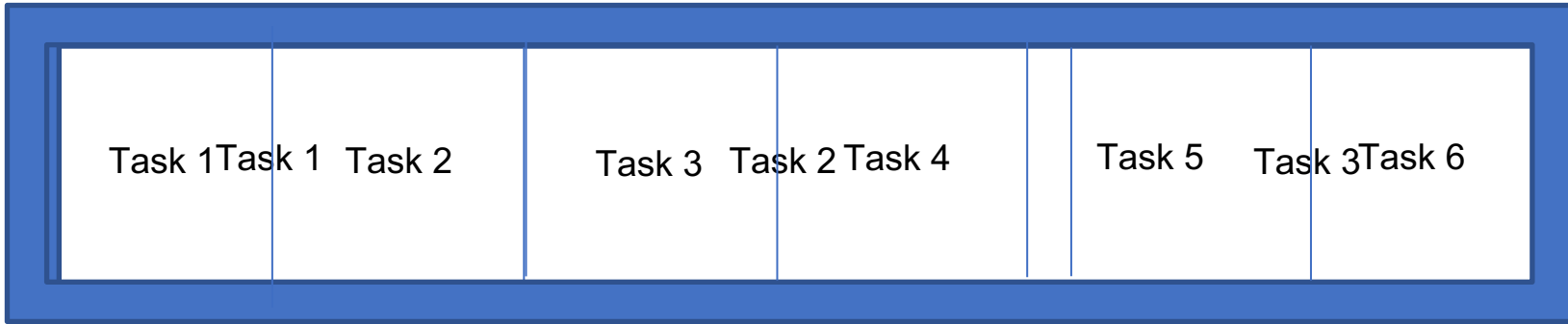
- `container_cpu_usage_seconds_total`



** Per Period*

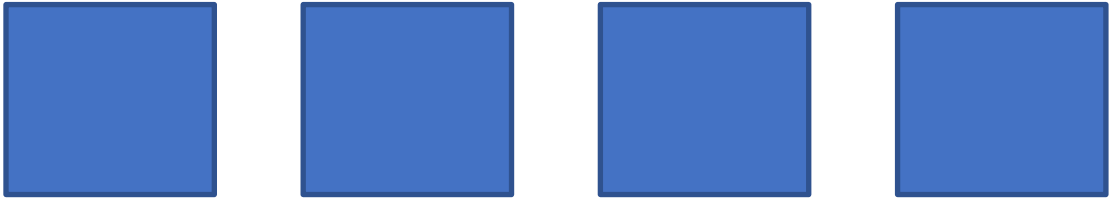
Threads

Container 1

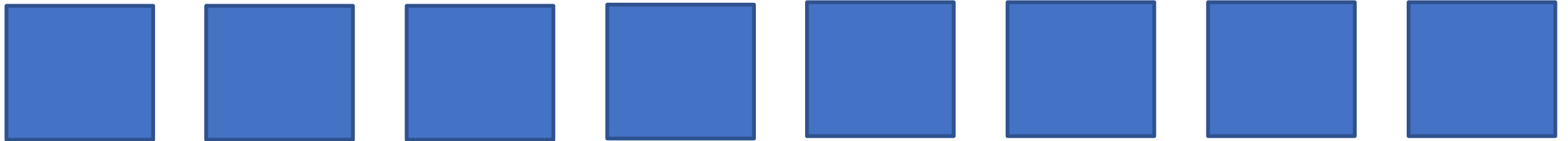


Swapping Cores

4 OS Threads



8 OS Threads



Throttles

Throttles - Bad Pods Limits	
pod full-test-27434725--1-b6qn7	99.8%
pod full-test-27434725--1-w55zh	99.8%
pod full-test-27434725--1-7zsdd	99.7%

 Query

 Transform

1

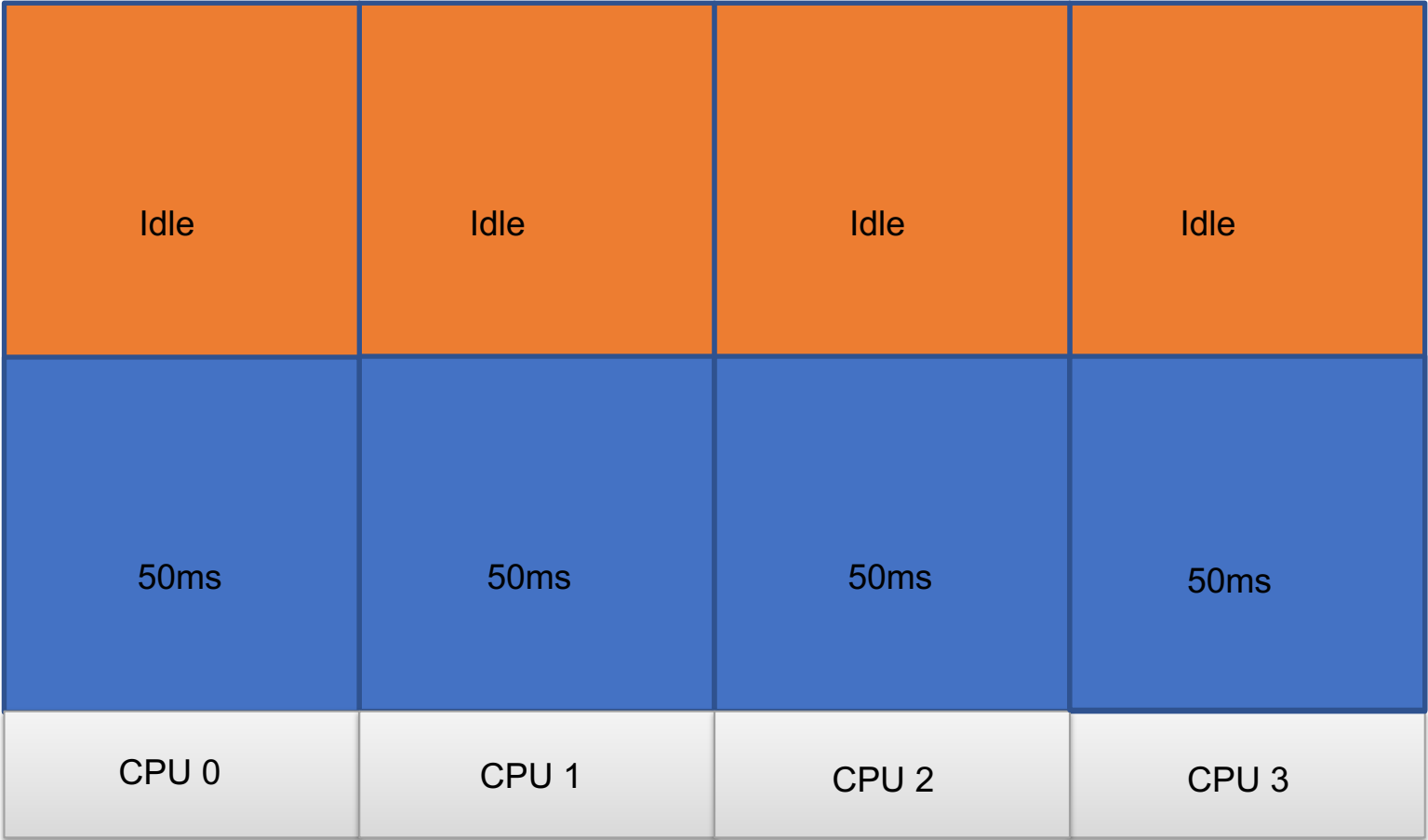
0

Double It!

Resources:
limits:
cpu: 2000m

or

200ms

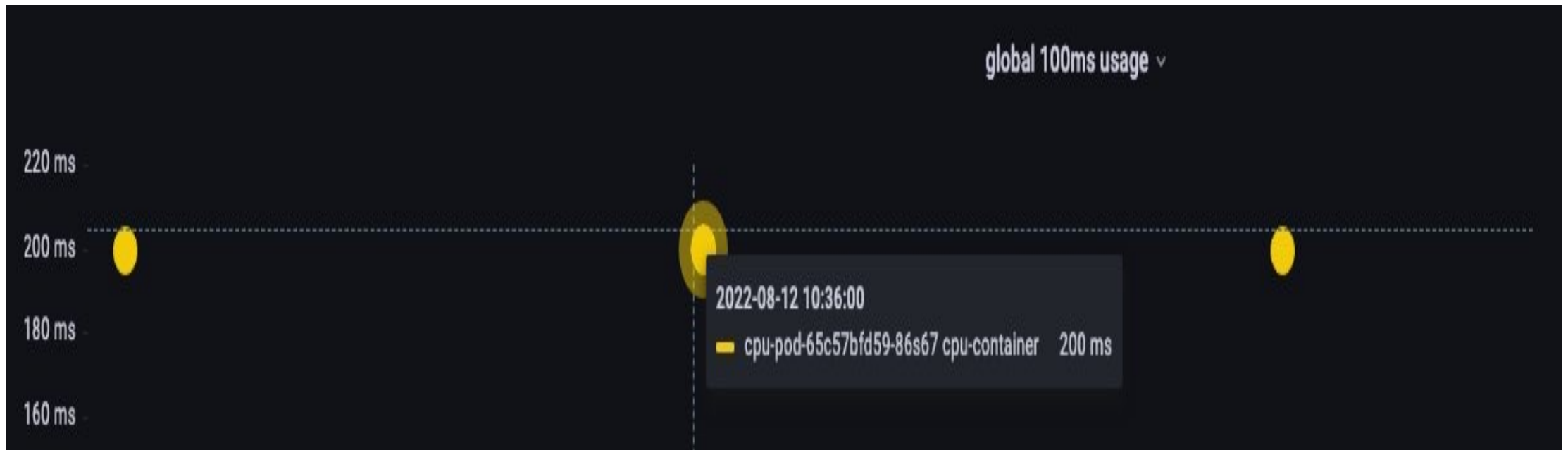


Real Time
100ms

Quota Time
200ms = 2 CPU

Thinking in Time

- `container_cpu_throttled_seconds_total`



* *CPU Usage Per Period*

Container
Global quota

Container
Quota 50ms

5ms slice

5ms slice

5ms slice

5ms slice

CPU 1

CPU 2

CPU 3

CPU 4

Thread 1

Thread 2

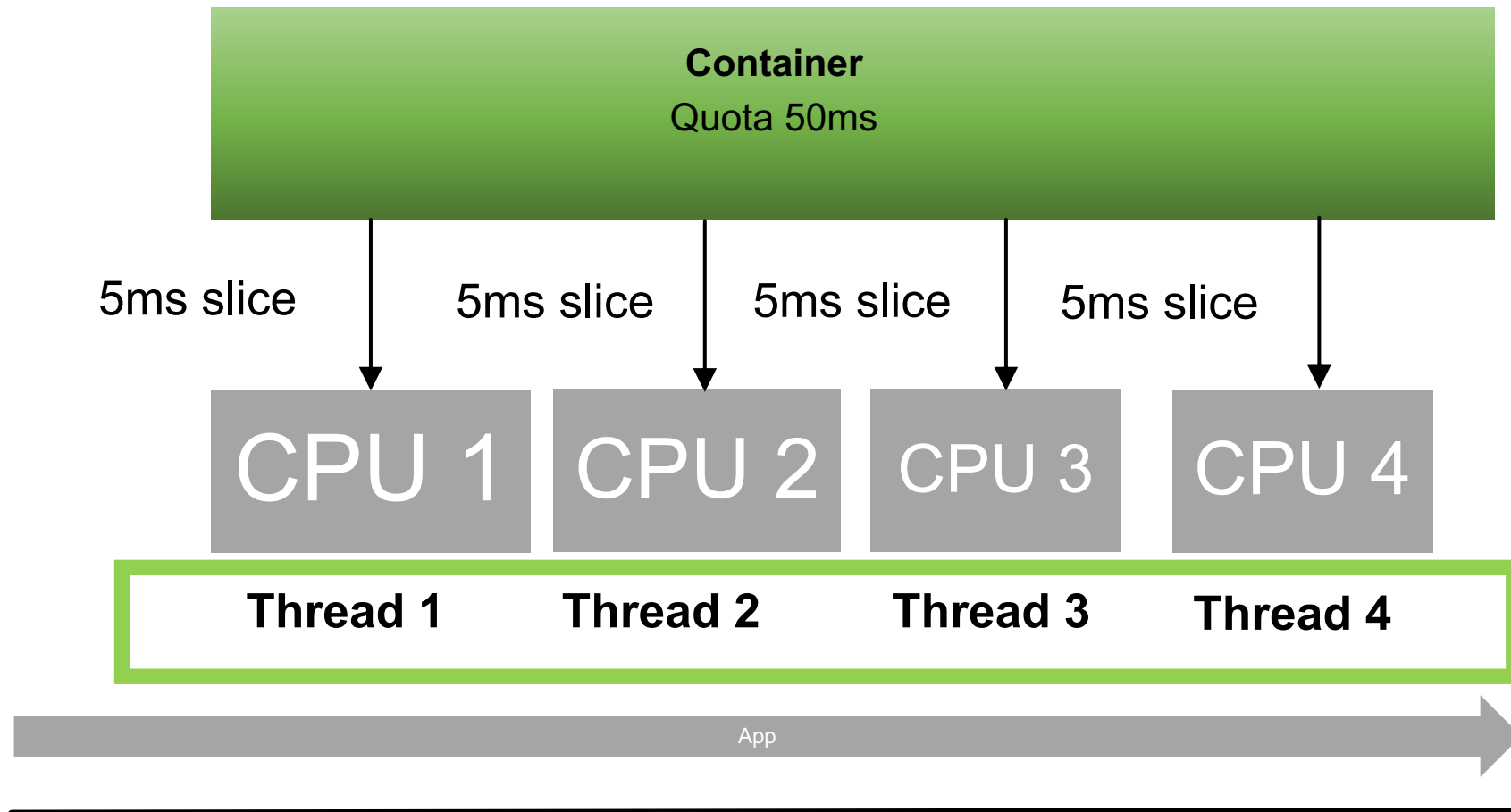
Thread 3

Thread 4

App

0ms

100ms





CONTAINERS

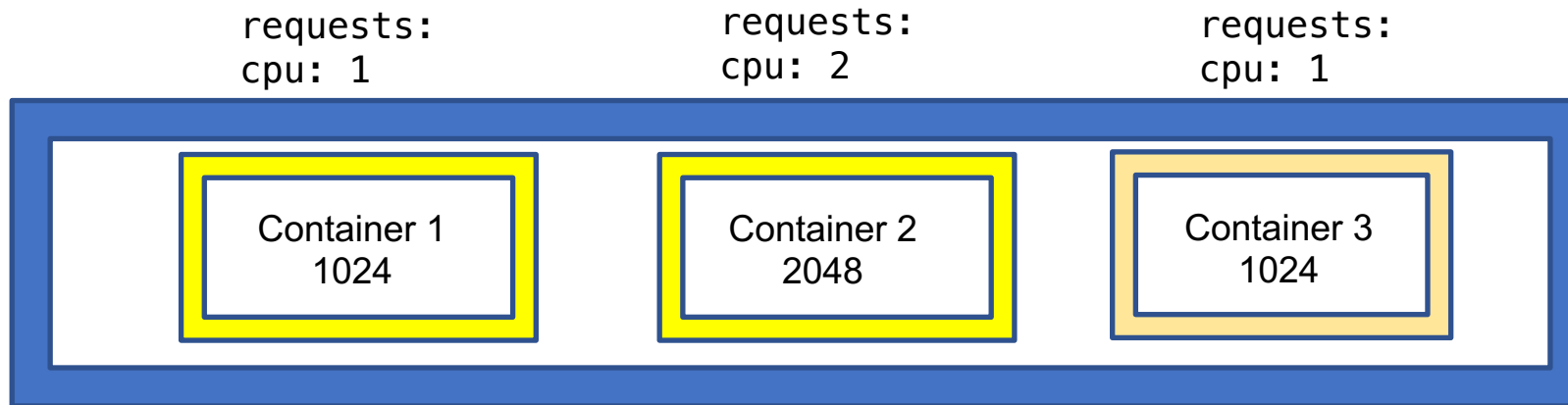
Using Prometheus to Avoid Disasters with Kubernetes CPU Limits

Read the blog post ›

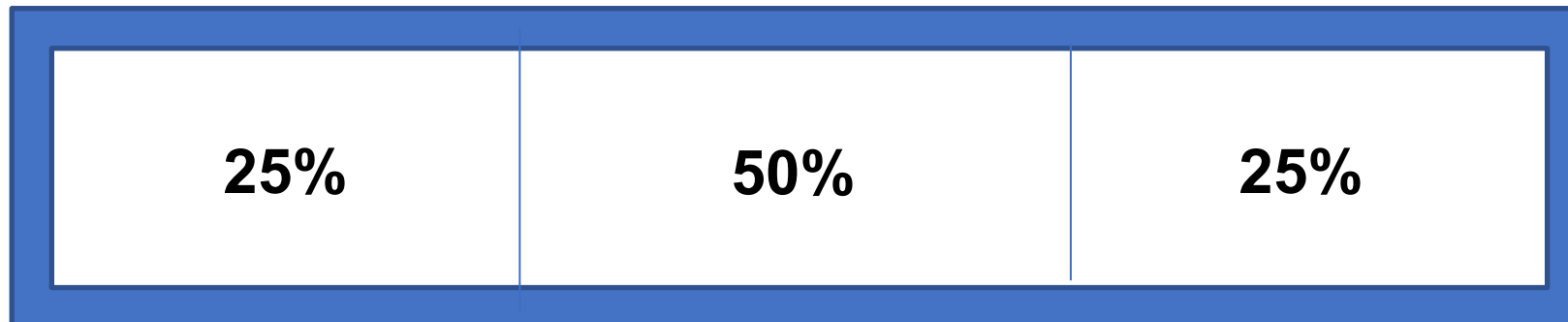


Static vs. Dynamic

Mixing Models

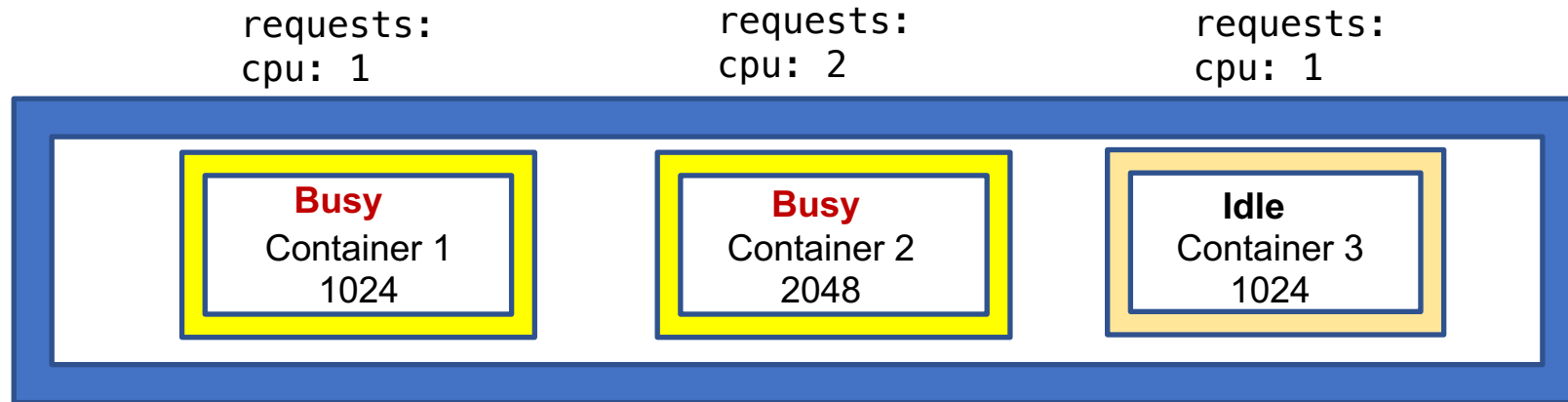


Linux Node

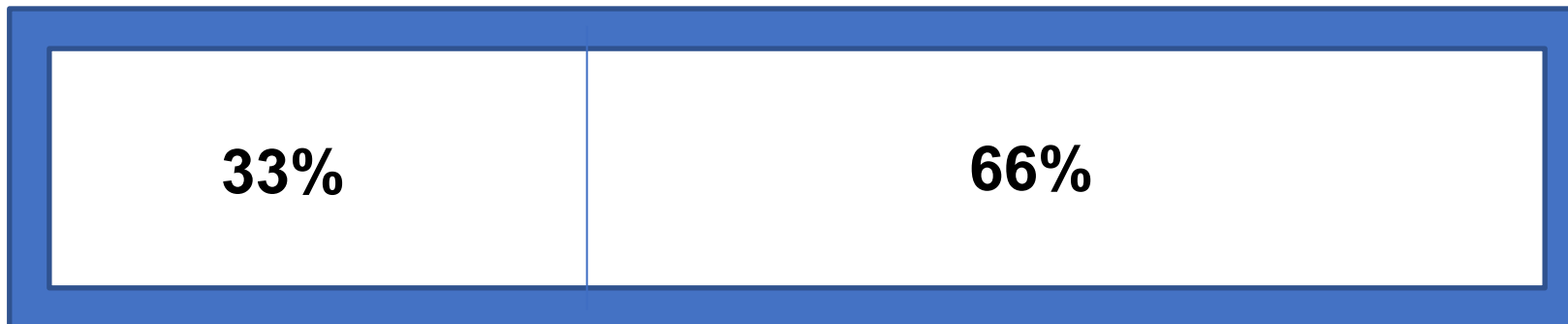


Dynamic Shares

$\text{CPU} * \text{Container Shares} / \text{Busy shares}$

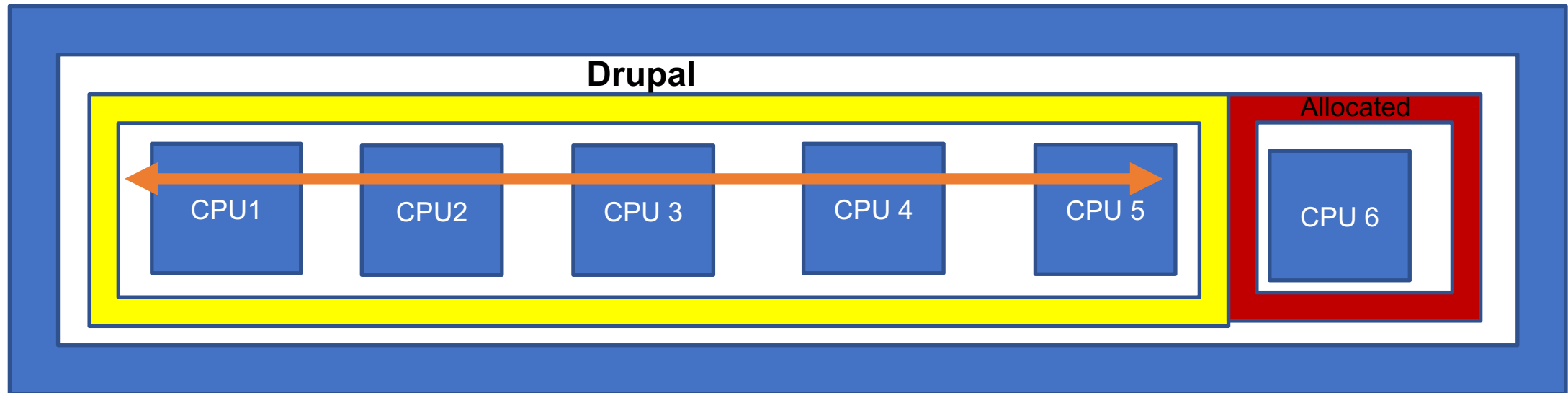


Linux Node

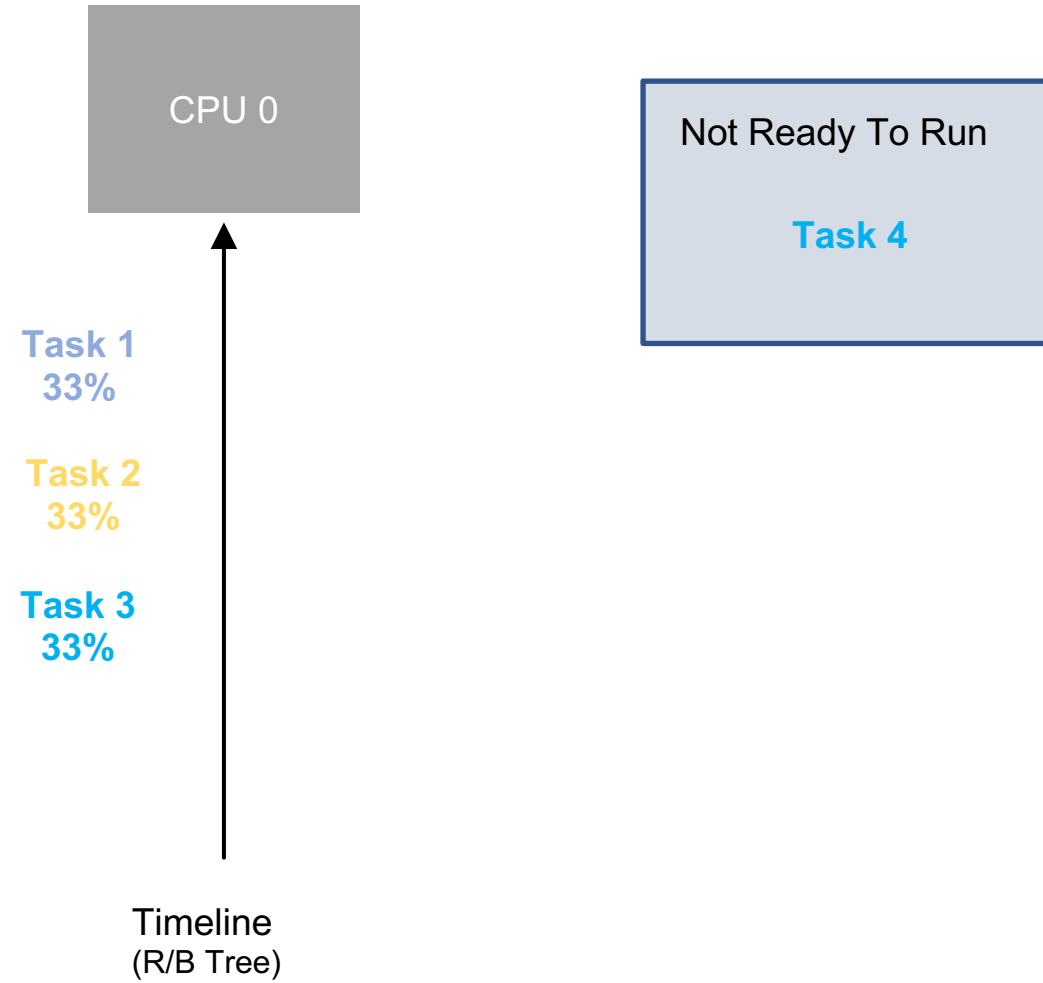


Only If It's Busy

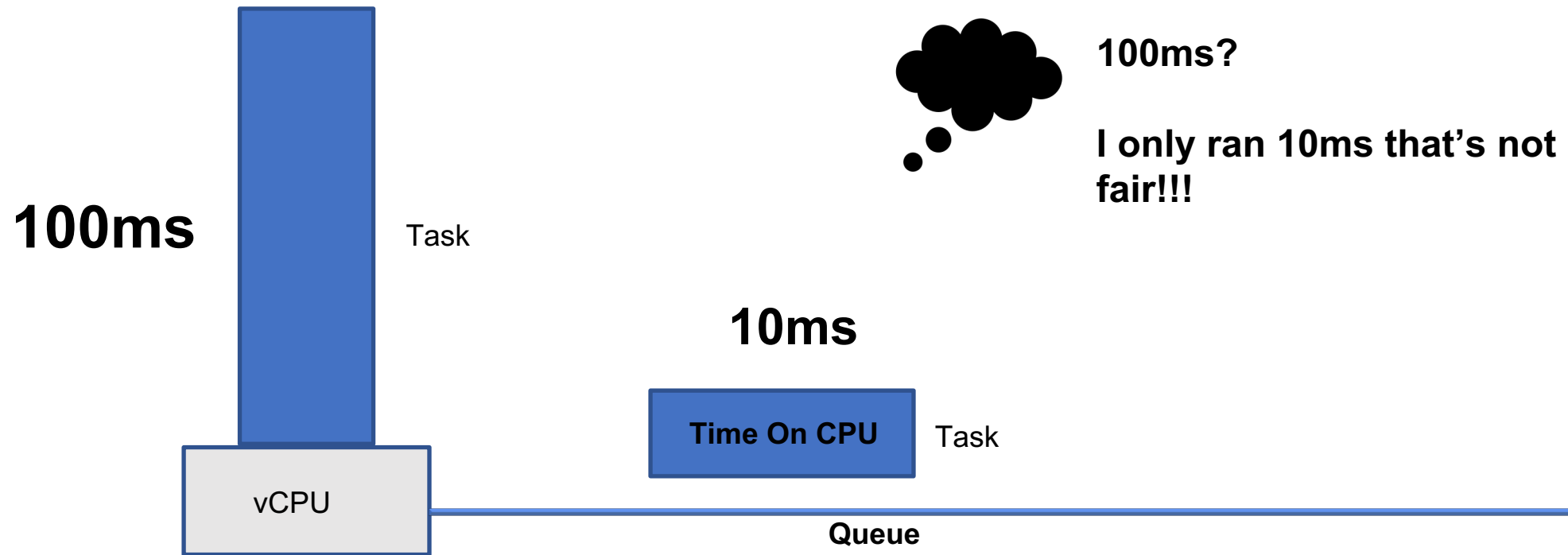
requests:
cpu: 1000m



Why?

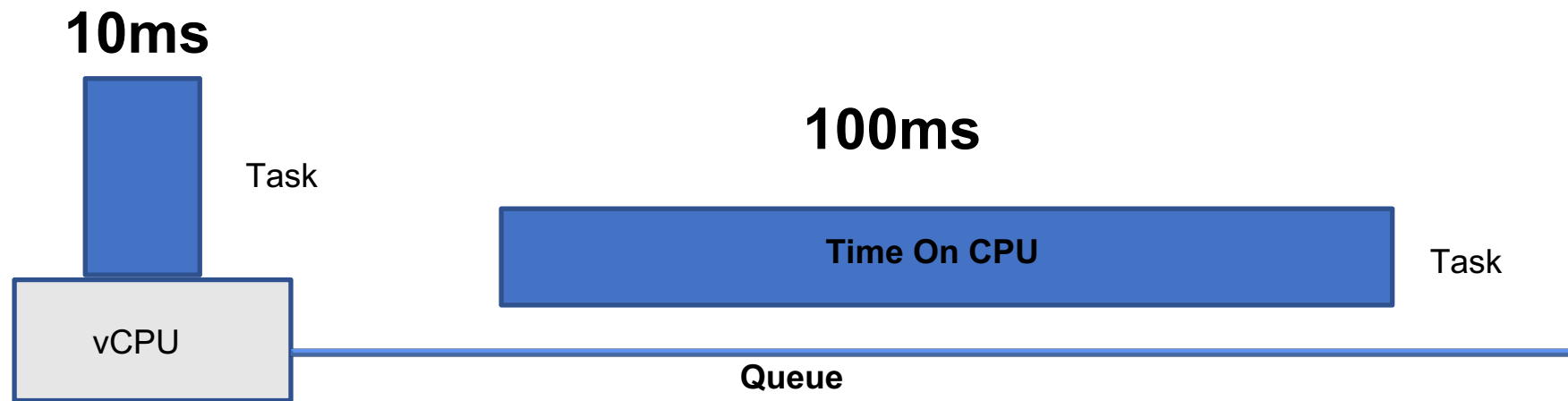


Linux is Fair



Linux is Fair

Lowest Time on CPU goes first!



Time Not Cores

TASK RUNTIME

100ms

requests:
cpu: 1

TASK RUNTIME

Runtime / 1024 shares

Shorter Runtime

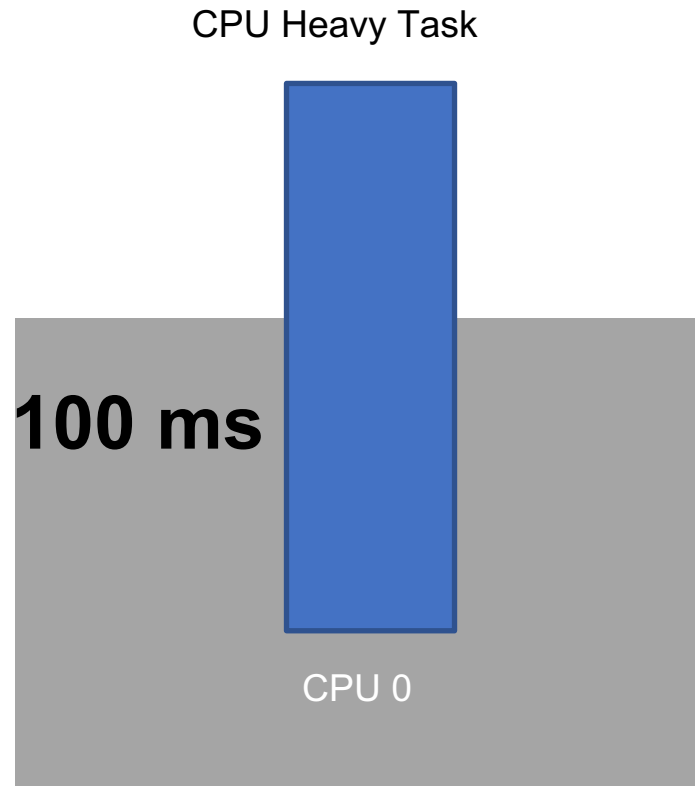
requests:
cpu: 2

TASK RUNTIME

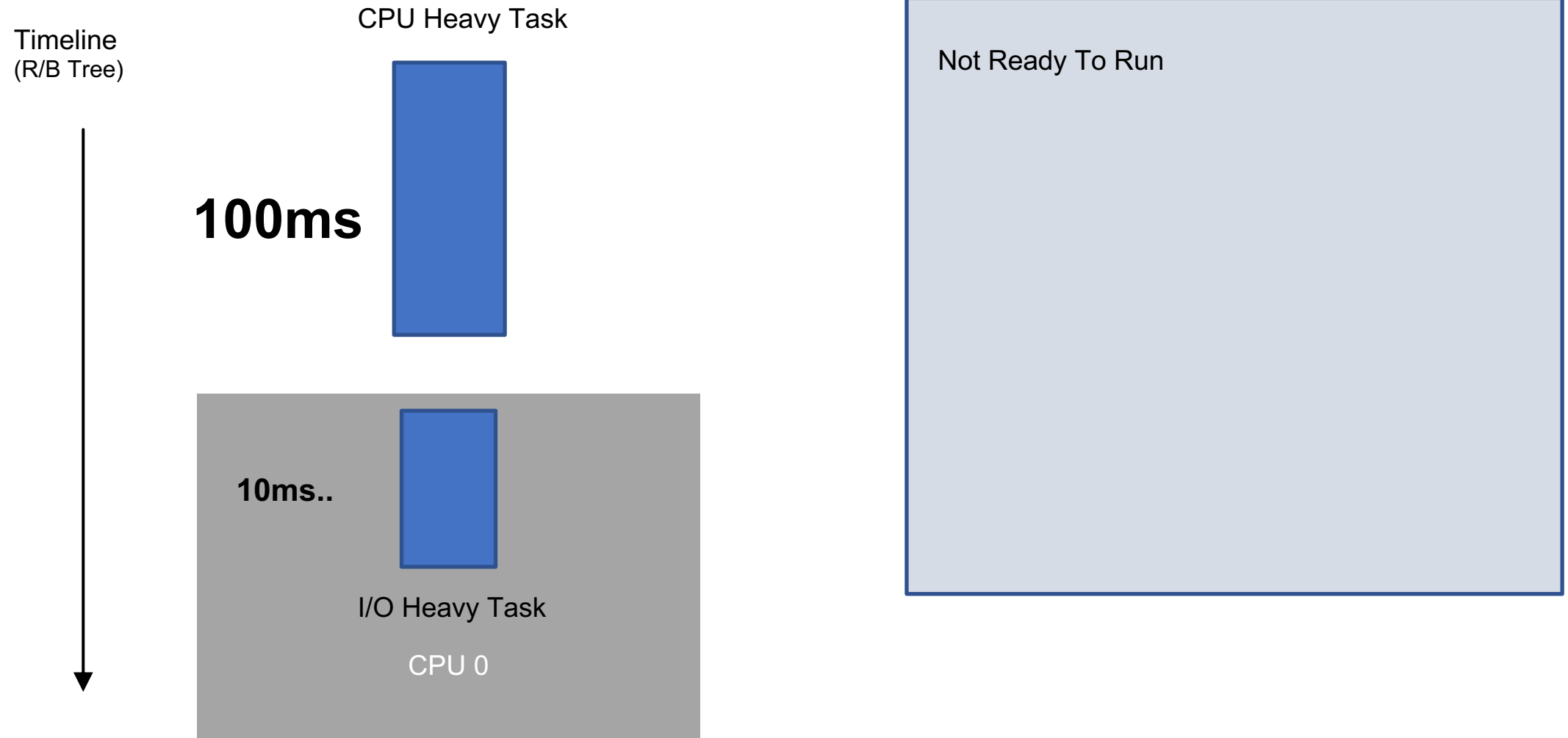
Runtime / 2048 shares

Shortest
Runtime

Least Amount of Time

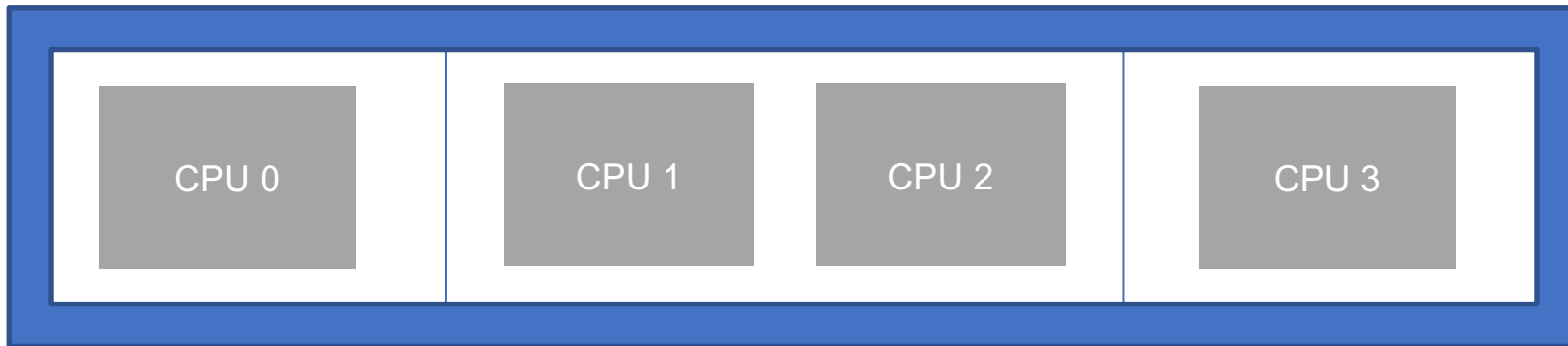
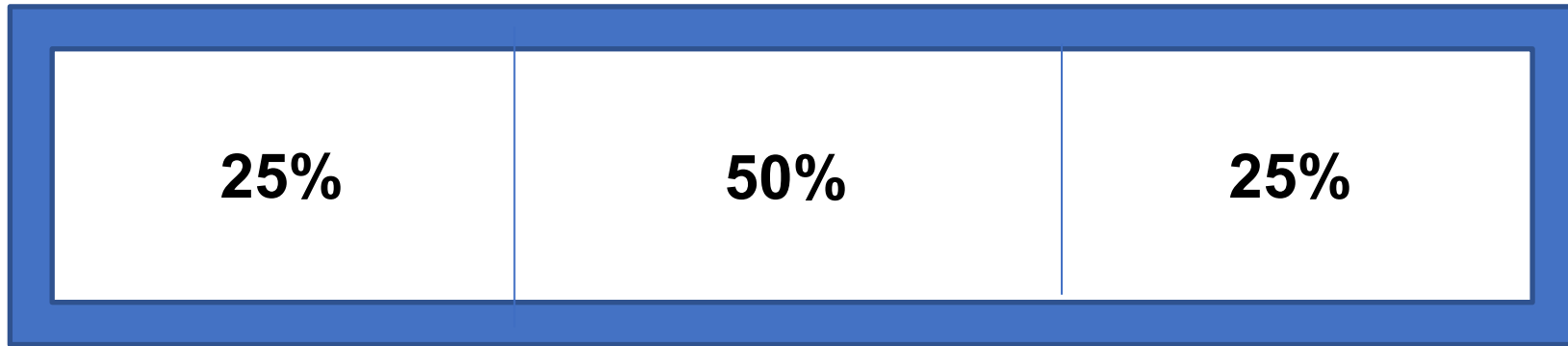


Don't Forget the I/O Bound

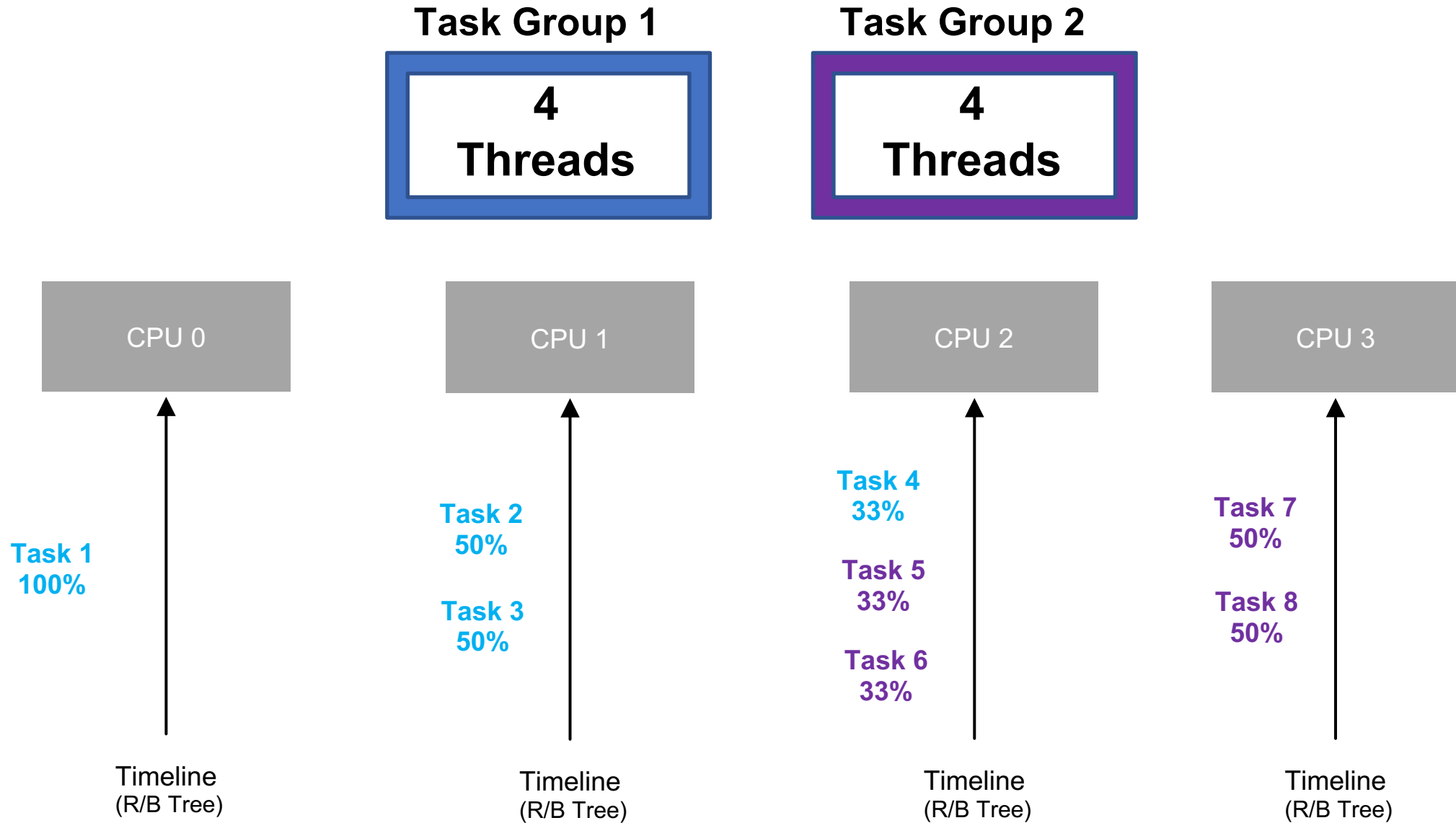


Not Cores

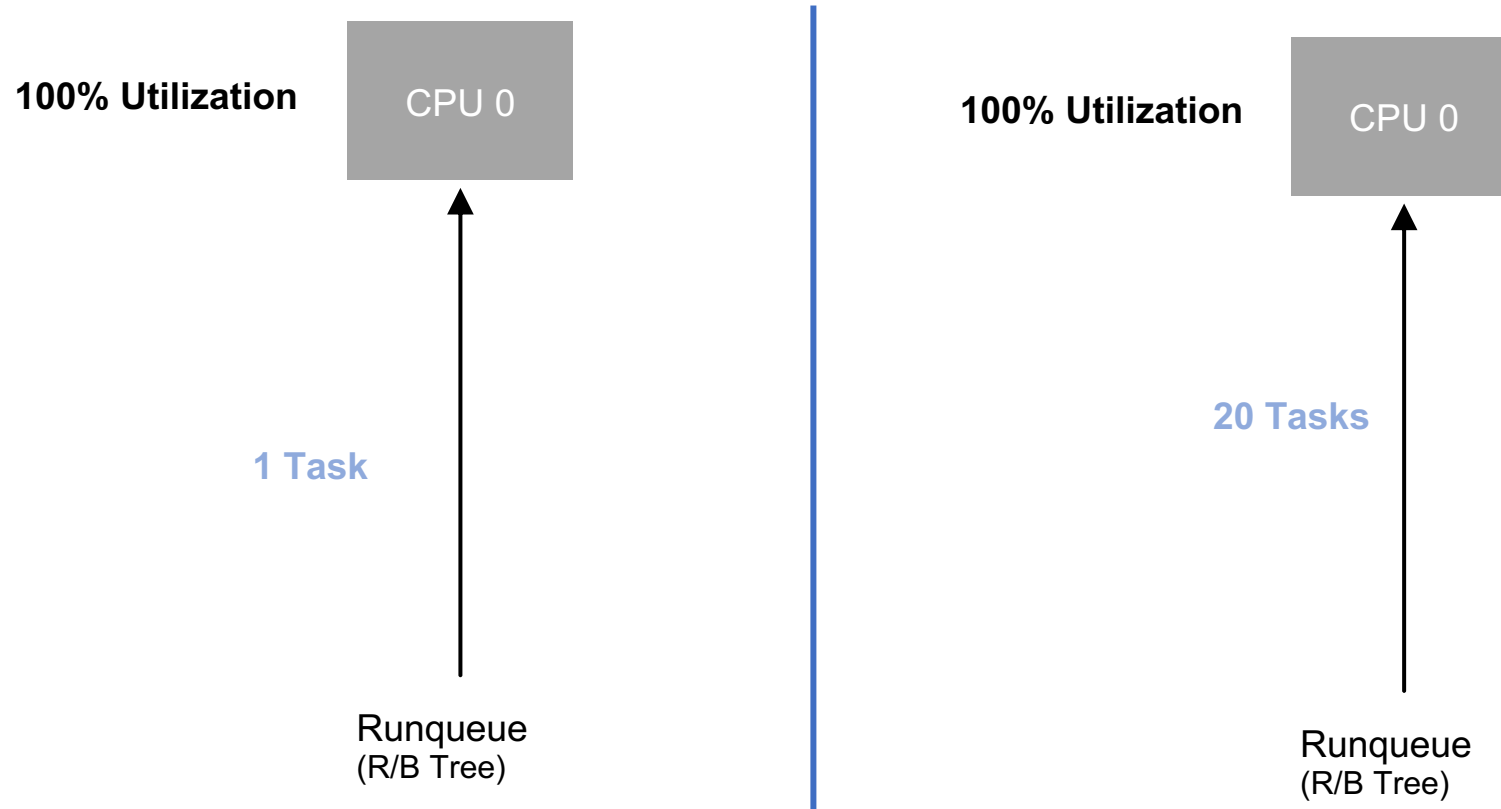
Linux Node



The New Picture

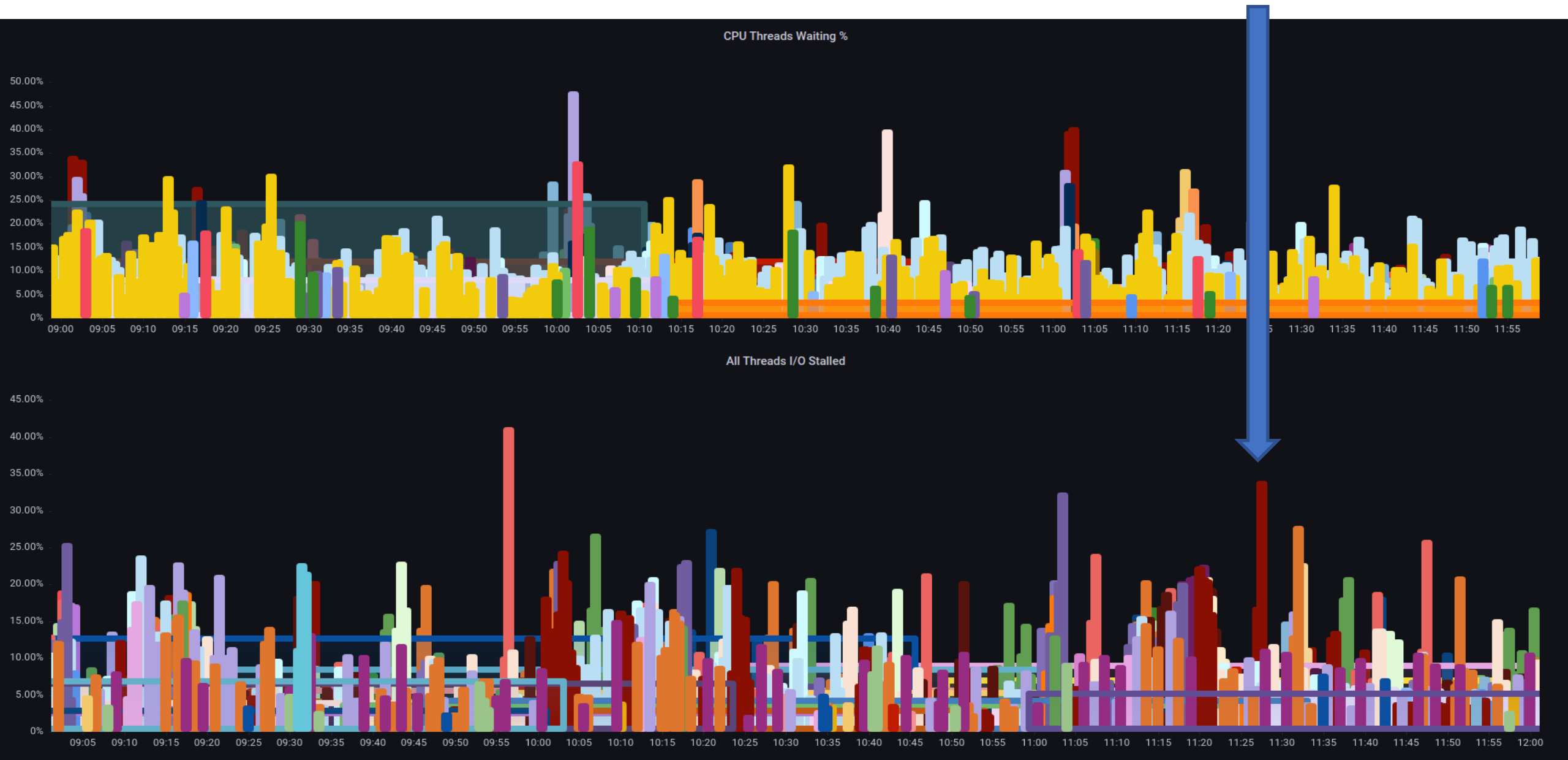


Utilization vs. Saturation



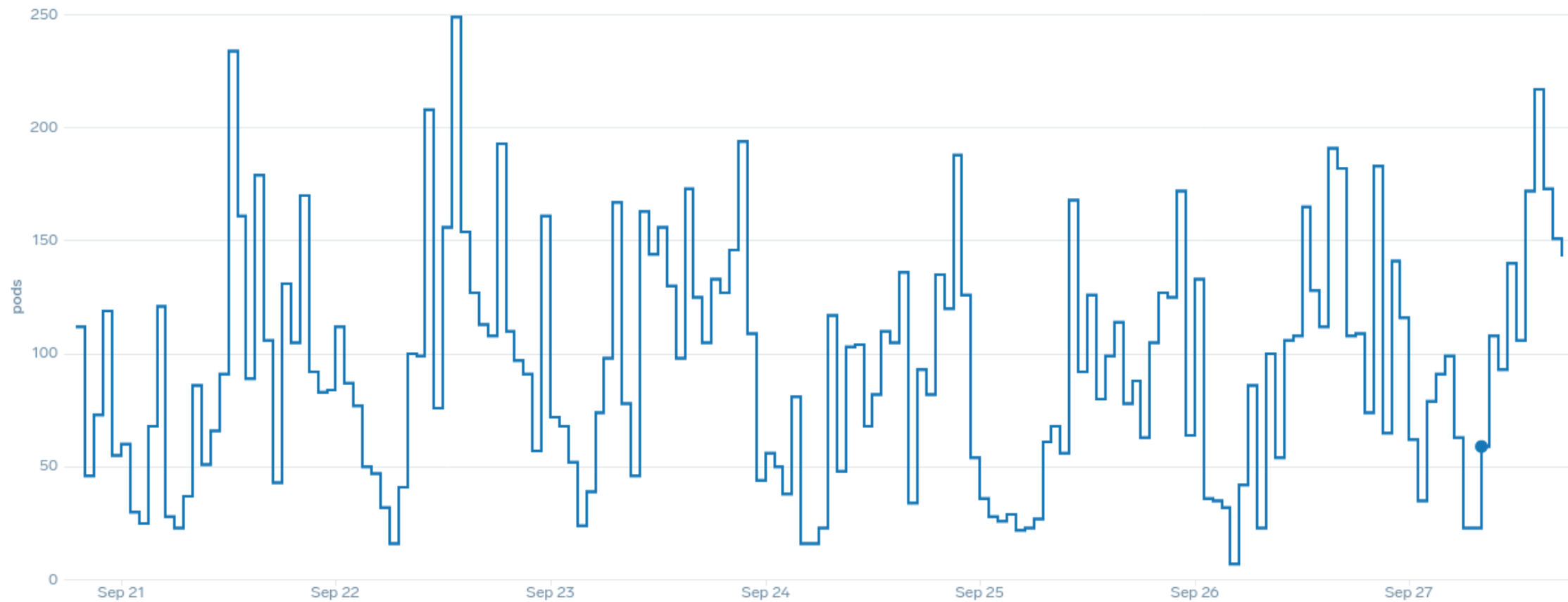
Pressure Stall Information

35% Stalled

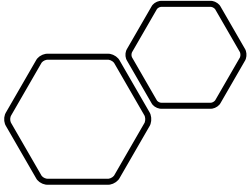




Number of pods

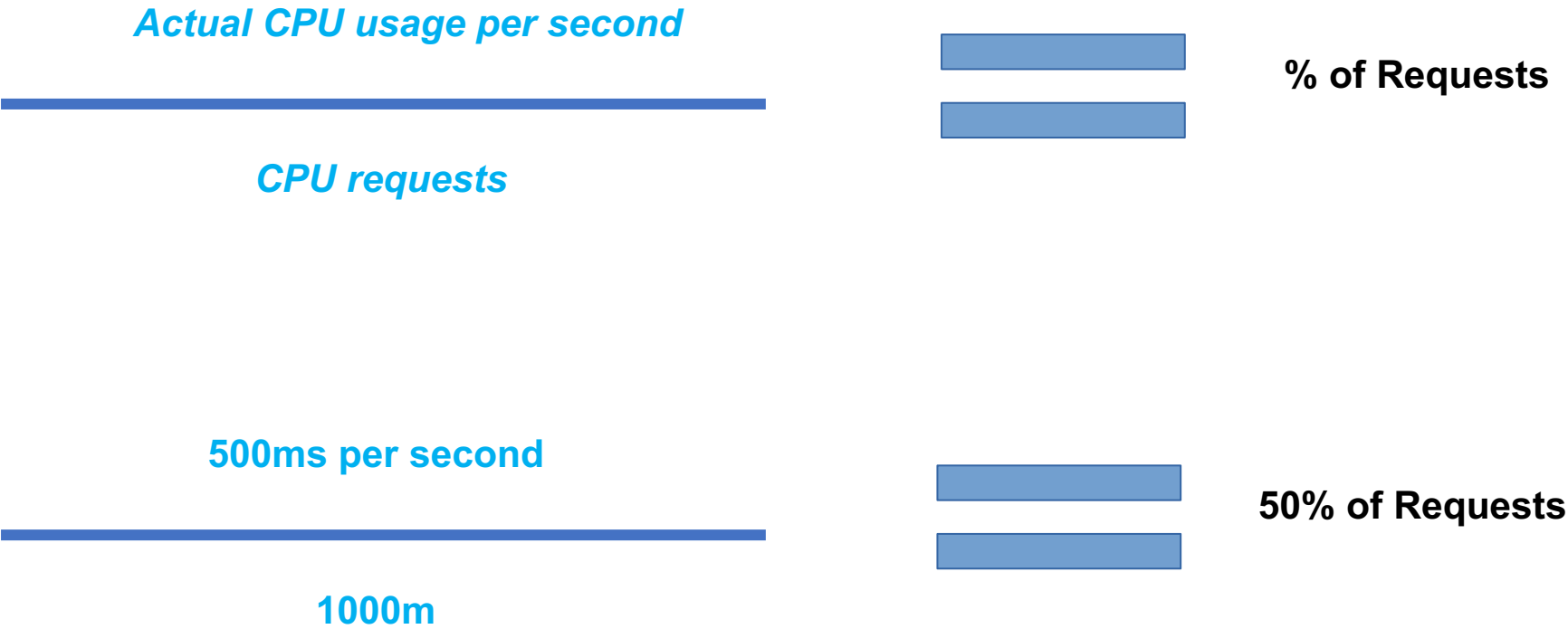


Unnecessary Churn



Percent of Requests?

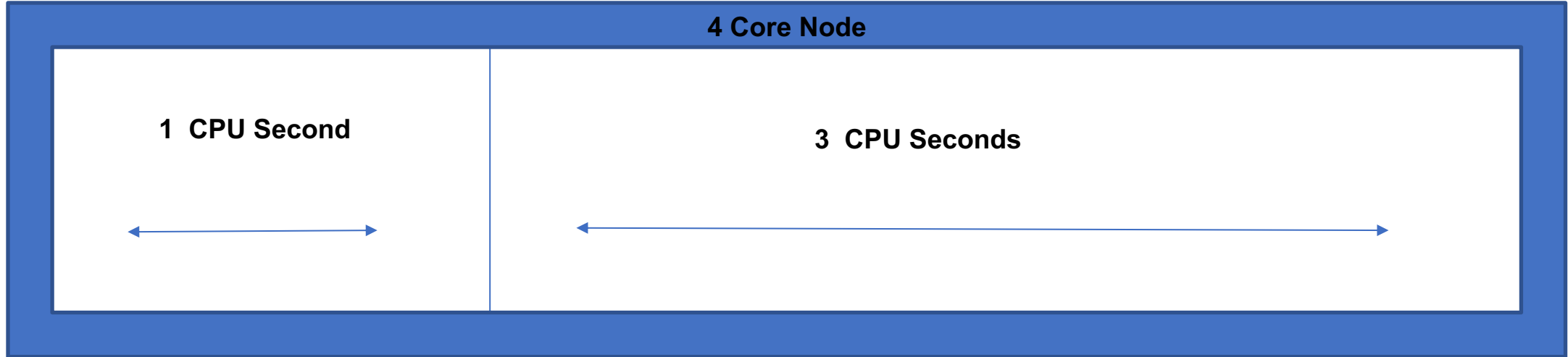
% Of Requests



100% of Request

requests:
cpu: 1

requests:
cpu: 3



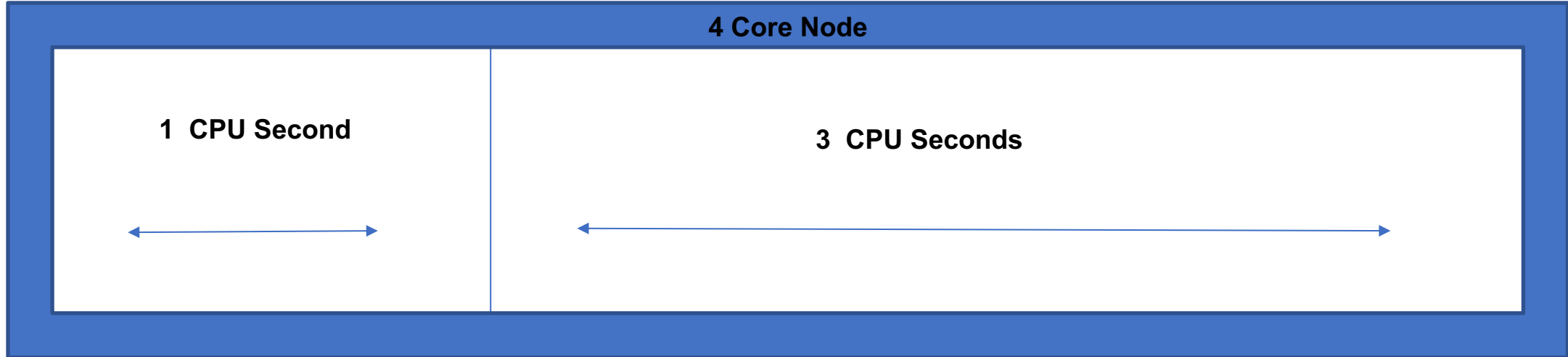
100% of Request

100% of Request

1,000% of Request

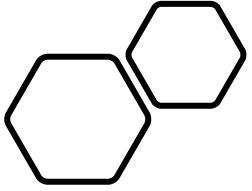
requests:
cpu: 100m

requests:
cpu: 300m



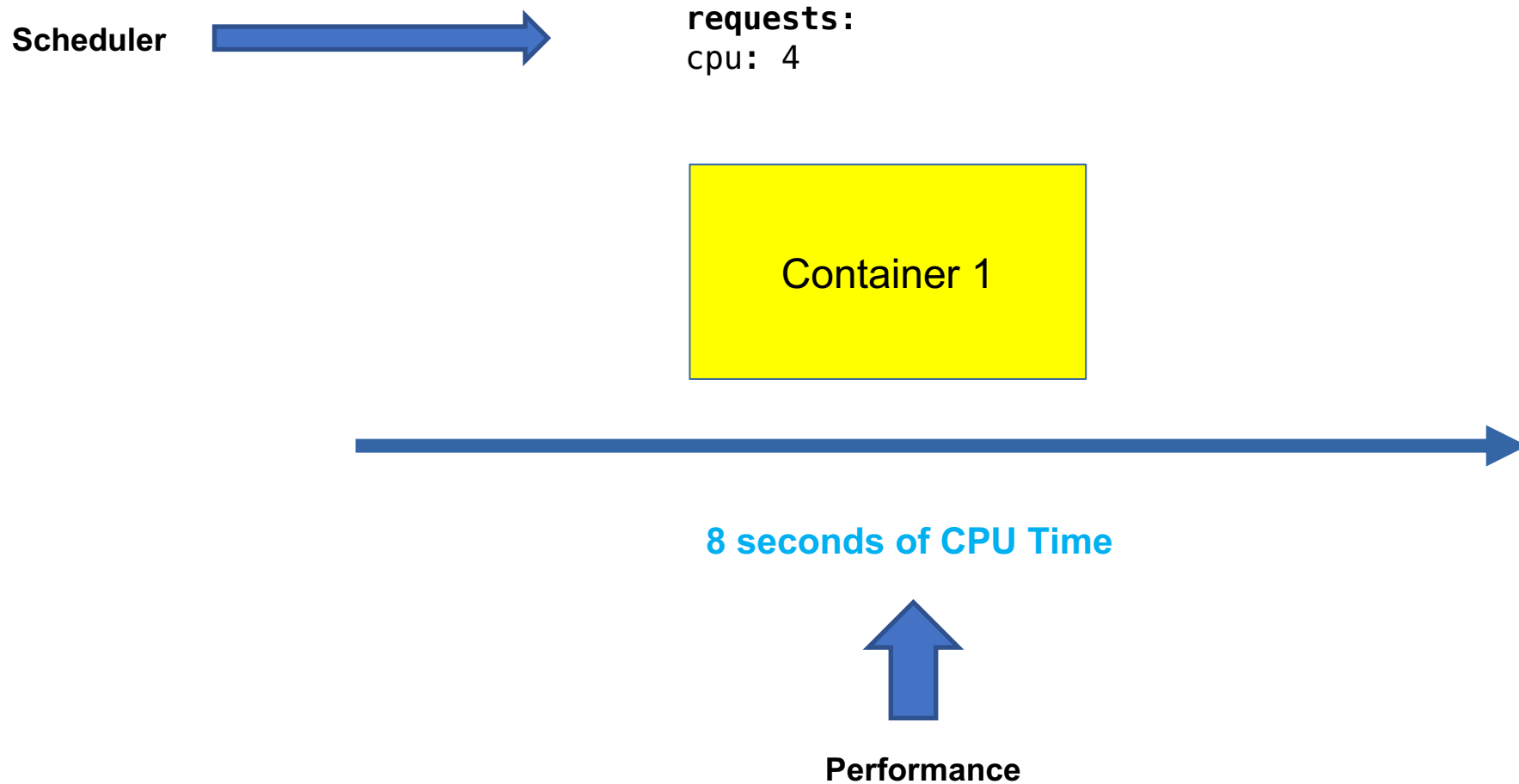
1,000% of Request

1,000% of Request

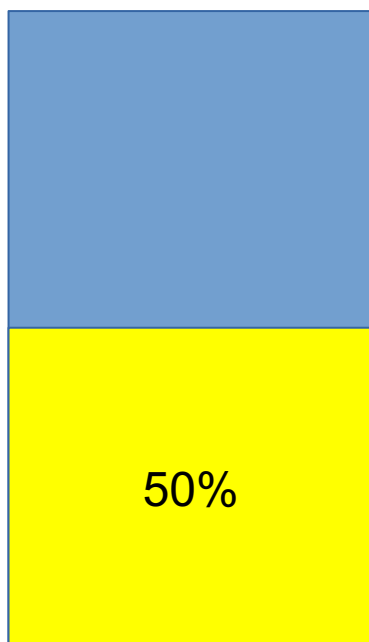
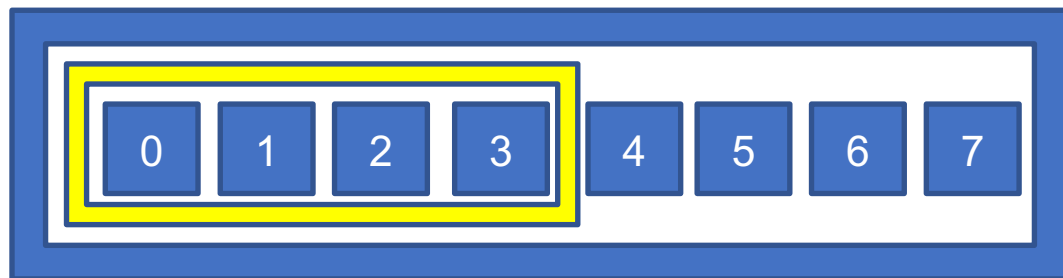


The Plot Thickens...

Scheduler vs. Performance



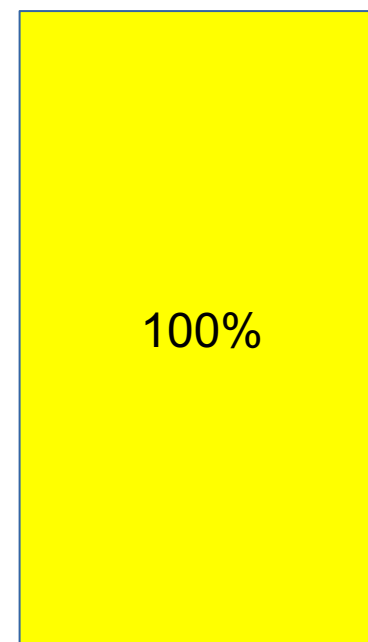
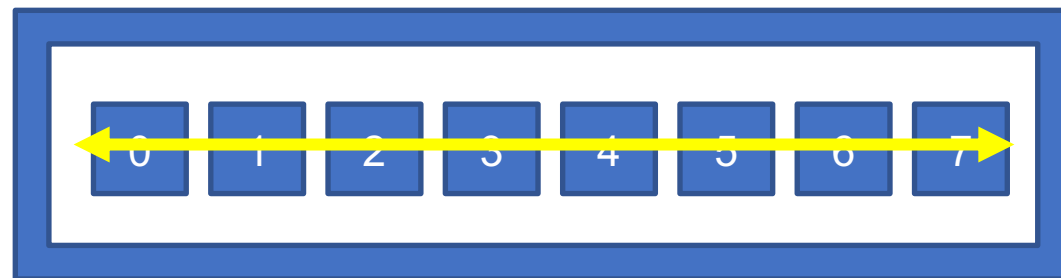
K8s Scheduler



Allocated

Linux

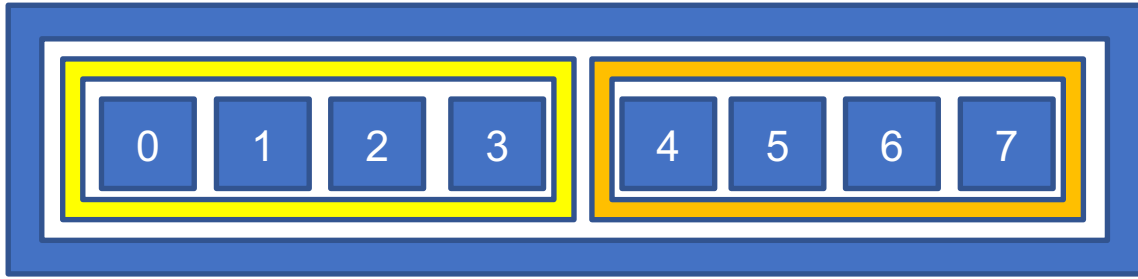
8 seconds of CPU Time



Real Time
1s

Utilized

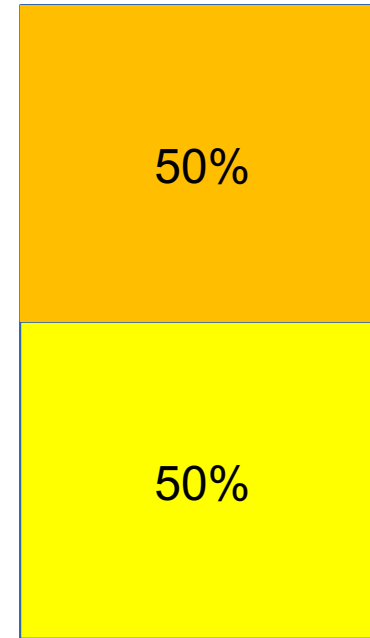
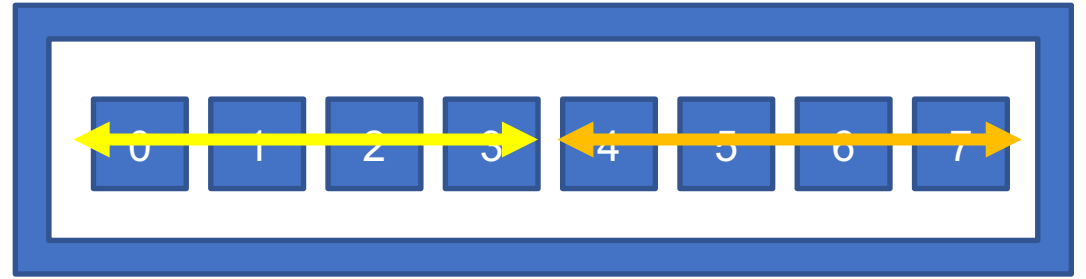
K8s Scheduler



Allocated

Linux

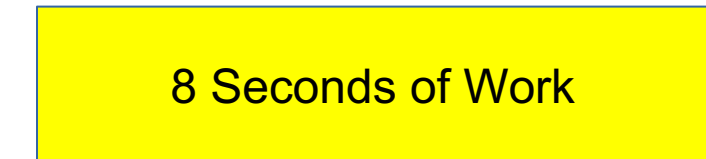
4 seconds of CPU Time | 4 seconds of CPU Time



Real Time
1s

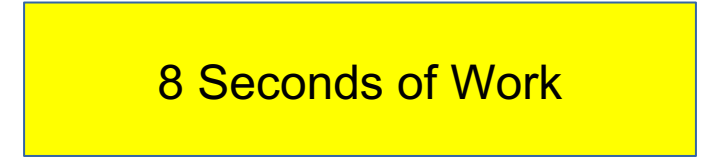
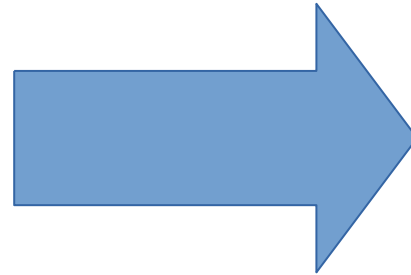
Utilized

Did we help?



1 second of Real Time

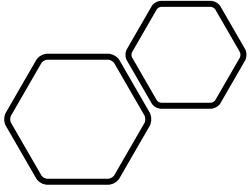
200% of requests



2 seconds of Real Time

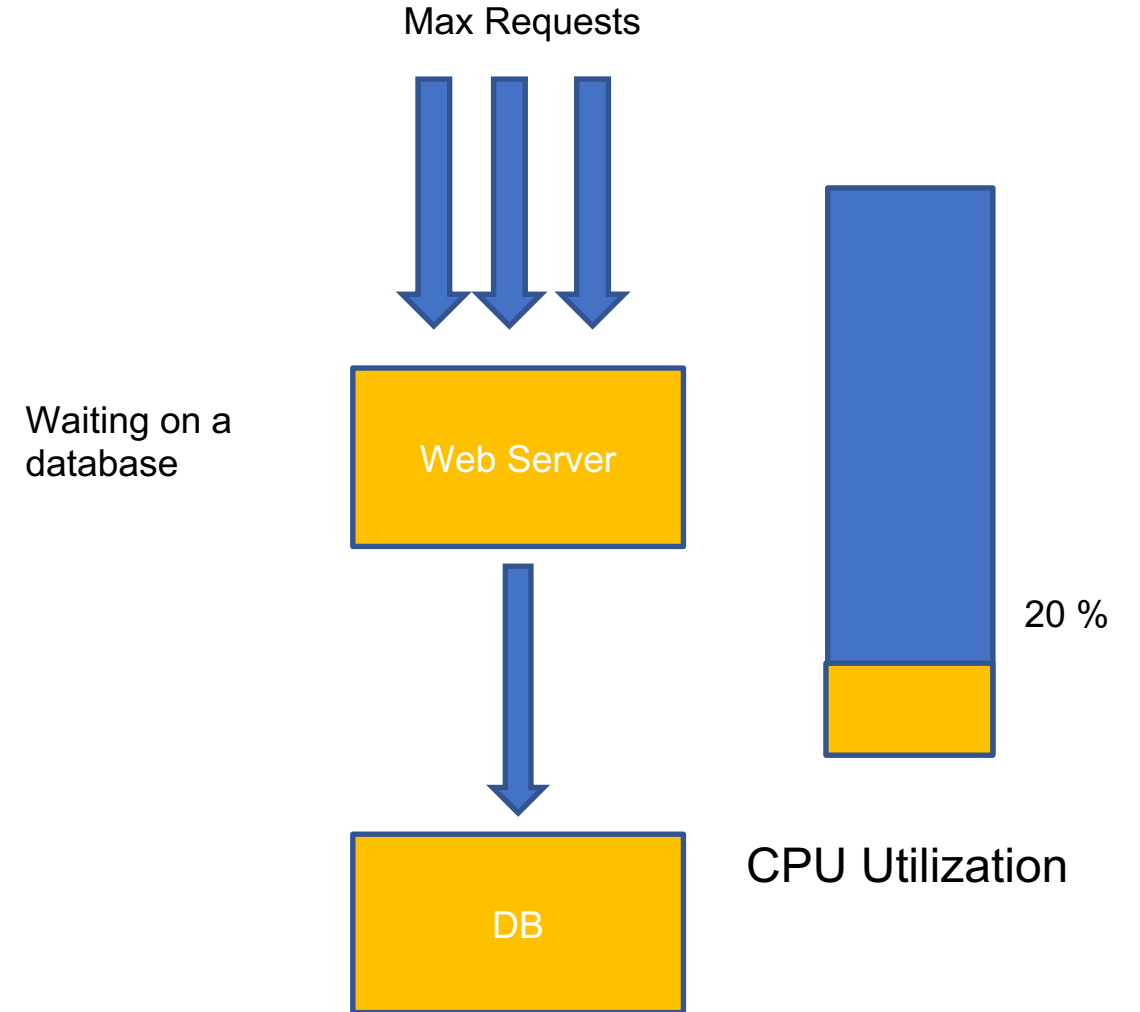
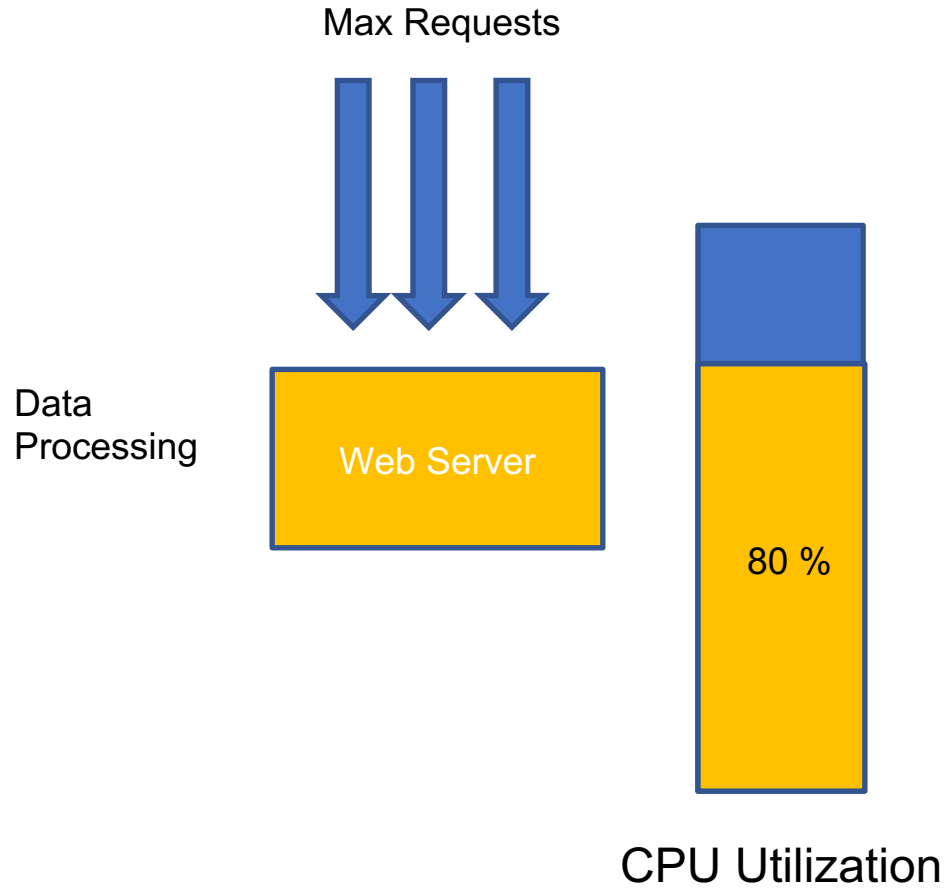
100% of requests



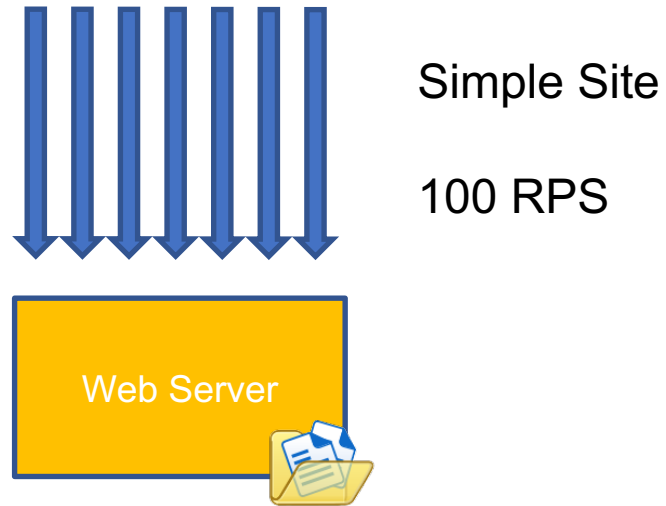


Is CPU The Right Metric At All?

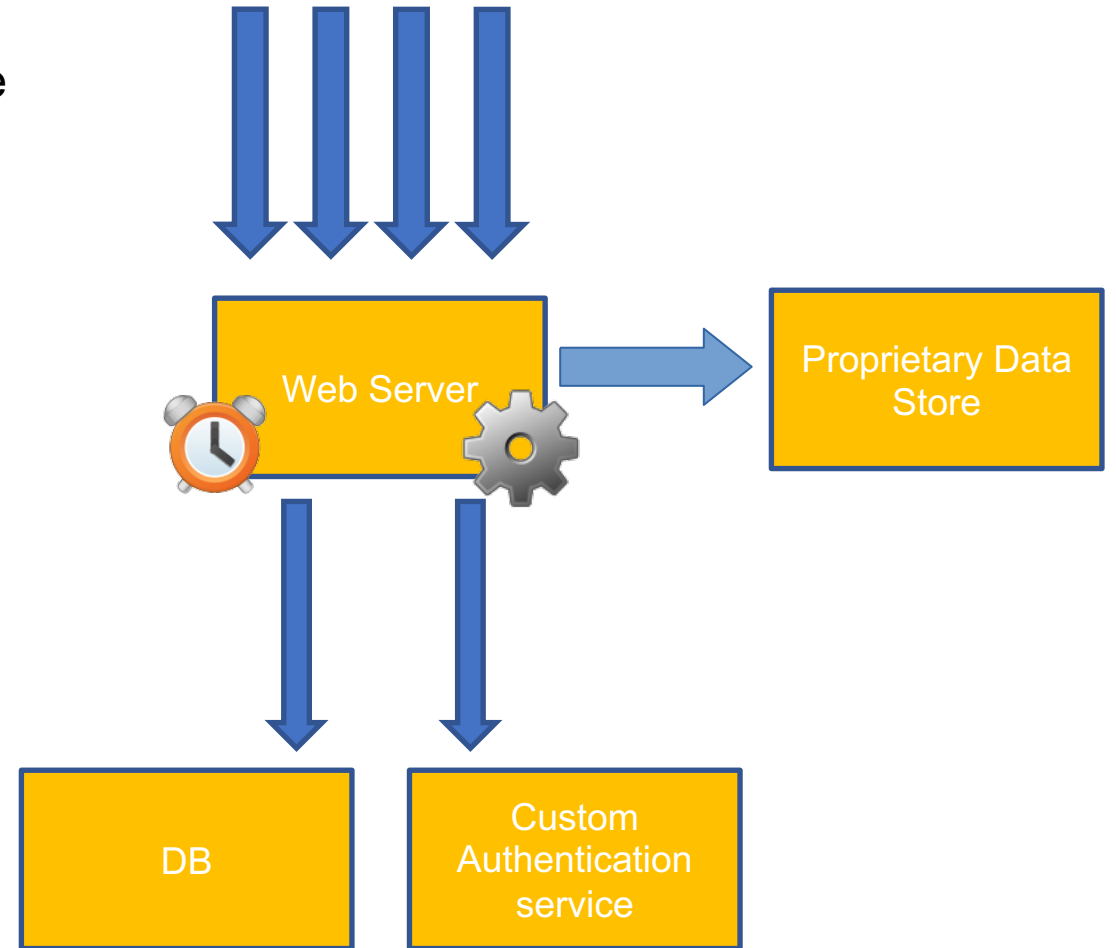
Saturated

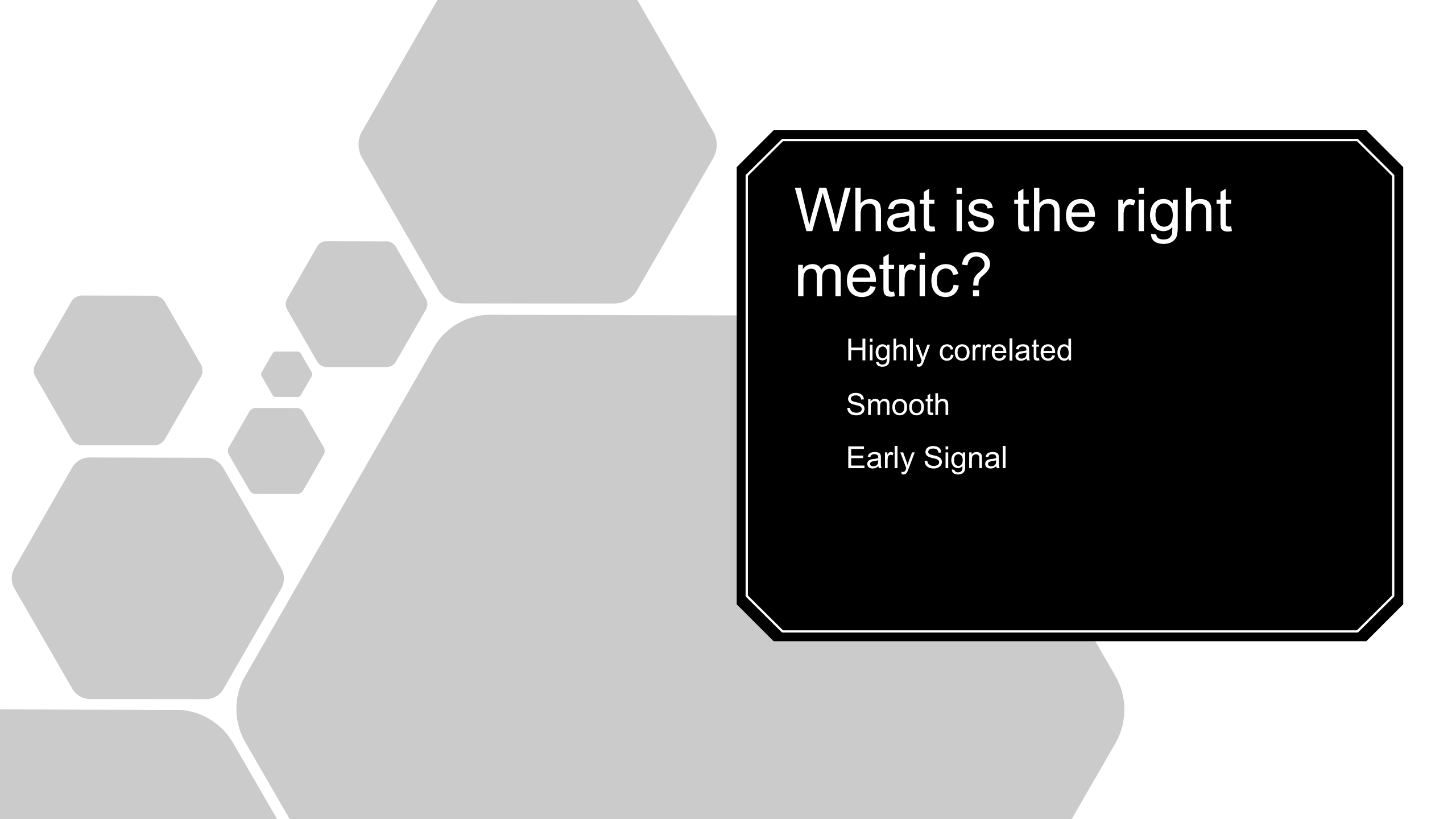


Not All Requests Created Equal



Complex site
4 RPS





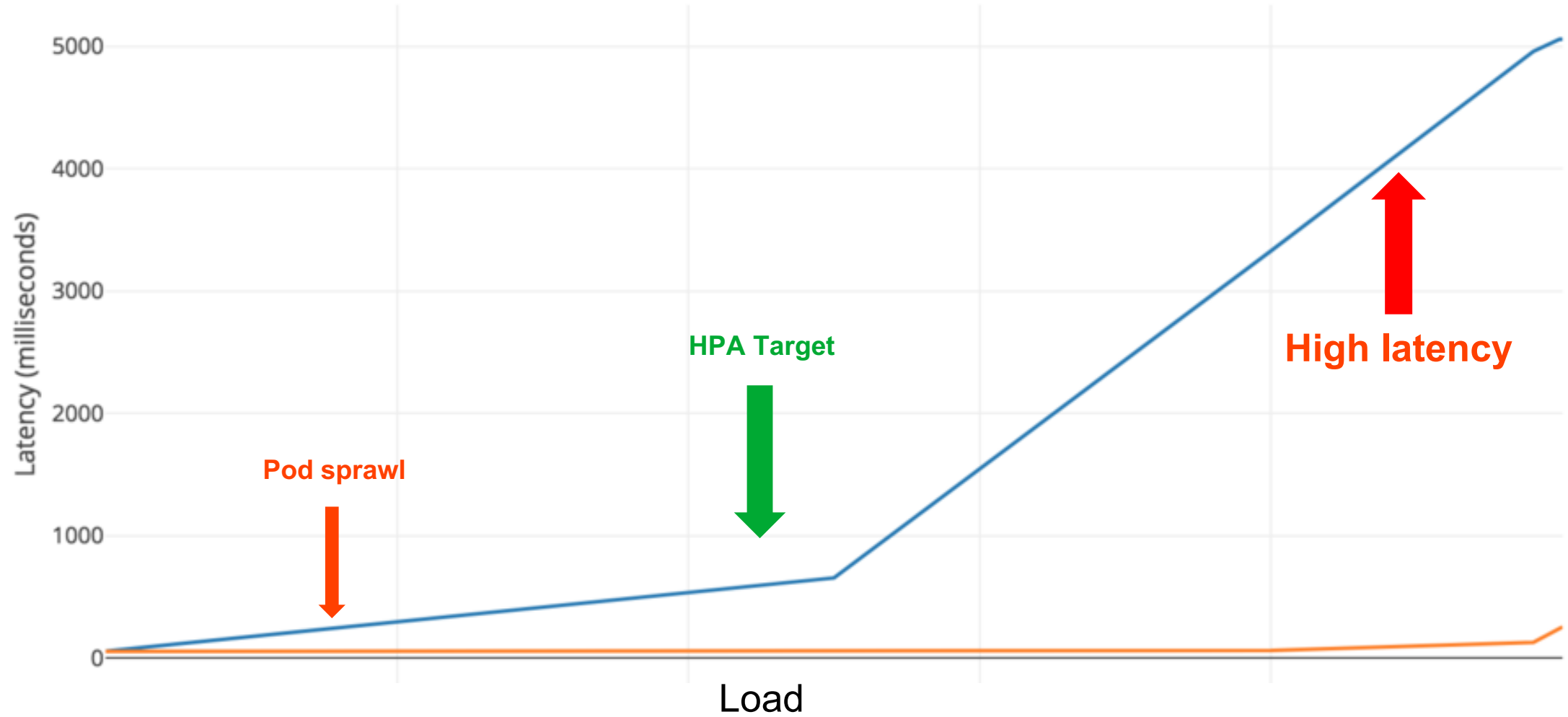
What is the right metric?

Highly correlated

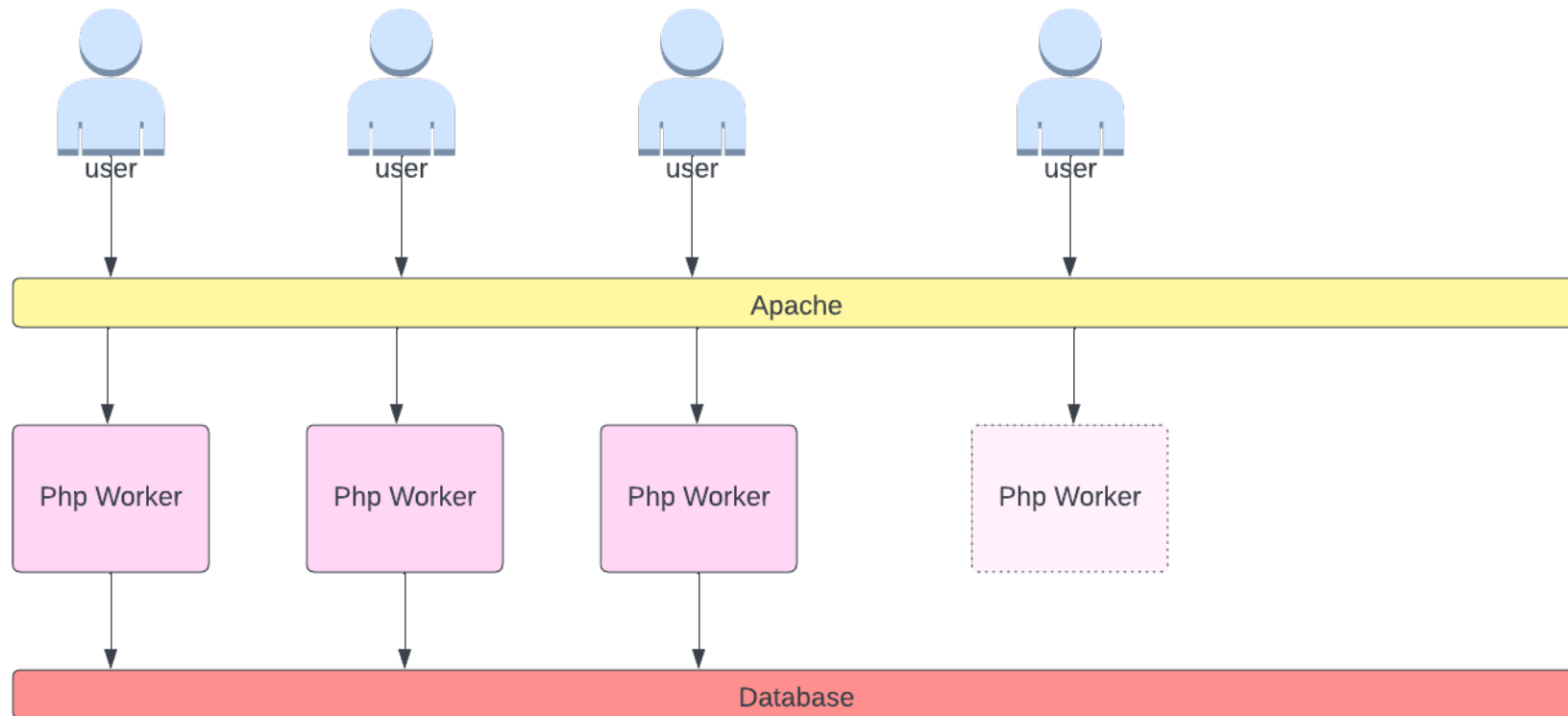
Smooth

Early Signal

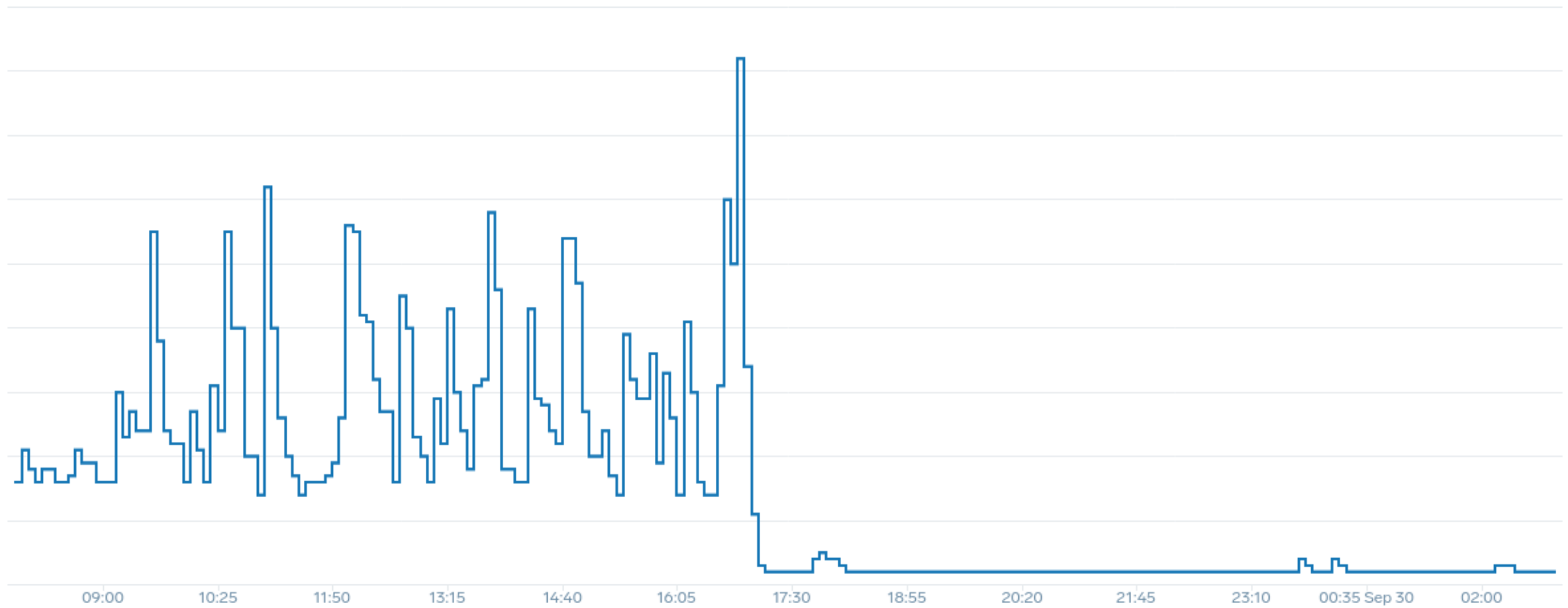
Workload Rightsizing



Active Threads



Success



Lessons

Focus on the Fundamentals

The Easy Answer is the Wrong One

Measure Everything

Don't Think in "Cores"

Don't Scale on CPU/Memory

Understand the Metrics

Session QR Codes will be sent via
email before the event

Please scan the QR Code above to
leave feedback on this session