

# Throttling: New Developments in

## Application Performance with CPU Limits

Dave Chiluk, Linux Platform Software Engineer

Welcome. Here's what we'll cover today.

---

- + How container CPU constraints work
- + Reproducing the throttling problem
- + The root cause
- + Solutions and workarounds



# Throttling: New Developments in Application Performance with CPU Limits

**Dave Chiluk**

Linux Platform Software Engineer, Indeed



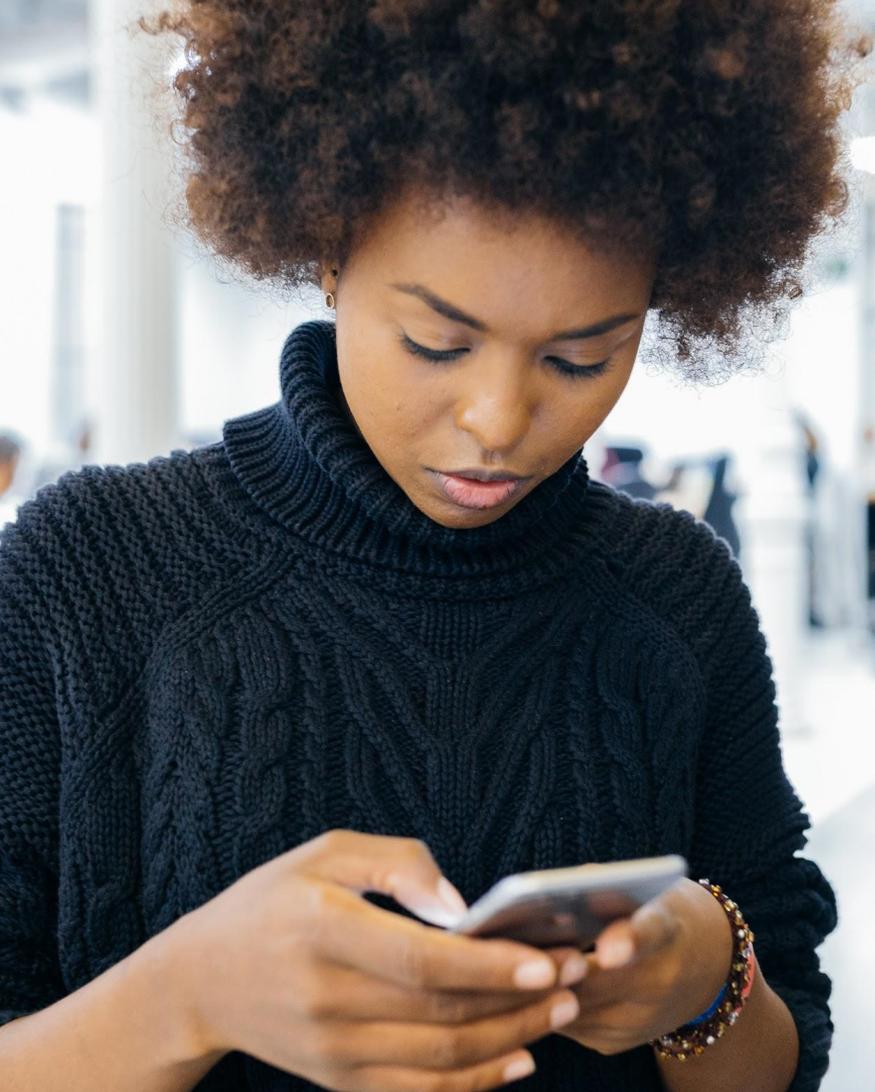
**Dave Chiluk**

Linux Platform Software Engineer

The logo for Indeed, featuring a white stylized 'i' with a curved line above it, followed by the word 'indeed' in a lowercase, rounded sans-serif font, all set against a solid blue background.

indeed

**We help  
people  
get  
jobs.**



**250M** unique monthly visitors<sup>1</sup>

**25M** jobs

**1M** employers with sponsored jobs

**10k** employees



<sup>1</sup> Source: Google Analytics, Unique Visitors, September 2018

# Roadmap



CPU Limit Basics



The Problem



Reproducing the  
Problem



Solution and  
Workarounds

# Roadmap



**CPU Limit Basics**



The Problem



Reproducing the  
Problem



Solution and  
Workarounds

**Who should care about CPU limits?**

EVERY CONTAINER  
ORCHESTRATOR  
ON THE PLANET



**kubernetes**



Apache  
**MESOS**<sup>™</sup>



**docker**

## Setting CPU Limits in Kubernetes

```
"resources": {  
  "limits": {  
    "cpu": "200m",  
  },
```

**Hard Limit**

---

```
  "requests": {  
    "cpu": "134m",  
  }  
}
```

**Soft Limit**

Located in pod definition @spec.containers[].resources.

# Soft Limits

$.requests = \text{Soft Limits} - \text{Cgroup cpu.shares}$

```
dchiluk@cando:/sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.shares  
1024
```

## Soft Limits

$$CpuTime = \frac{shares}{\sum shares} * NumCPUs$$

# Soft Limits

## Actual Usable CPU on 88 Core Machine



## Soft Limits

The Floor for Usable CPU

# Hard Limits

.limits = Hard Limits - Cgroup CFS Bandwidth Control

Containers are limited to using quota amount of CPU time in a period.

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.cfs_period_us  
100000  
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.cfs_quota_us  
10000
```

# THROTTLING



## Hard Limits

The Ceiling for Usable CPU

# Throttling metrics

Throttled time

nr\_periods

nr\_throttled

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.stat  
nr_periods 128  
nr_throttled 124  
throttled_time 5664985136
```

## Throttling metrics

### Throttled time

nr\_periods

nr\_throttled

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.stat  
nr_periods 128  
nr_throttled 124  
throttled_time 5664985136 ←
```

**Throttled time** is the sum total time a thread in a cgroup was throttled

## Throttling metrics

Throttled time

**nr\_periods**

nr\_throttled

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.stat
nr_periods 128 ←
nr_throttled 124
throttled_time 5664985136
```

Throttled time is the sum total time a thread in a cgroup was throttled

**nr\_periods** is the number of periods the application was running

## Throttling metrics

Throttled time

nr\_periods

nr\_throttled

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.stat
nr_periods 128
nr_throttled 124 ←
throttled_time 5664985136
```

Throttled time is the sum total time a thread in a cgroup was throttled

nr\_periods is the number of periods the application was running

**nr\_throttled** is number of those periods the application was throttled

## Throttling metrics

```
dchiluk@cando: /sys/fs/cgroup/cpu,cpuacct/user.slice/fibtest$ cat cpu.stat
nr_periods 128
nr_throttled 124
throttled_time 5664985136
```

Throttled time

Throttled time is the sum total time a thread in a cgroup was throttled

nr\_periods

nr\_periods is the number of periods the application was running

nr\_throttled

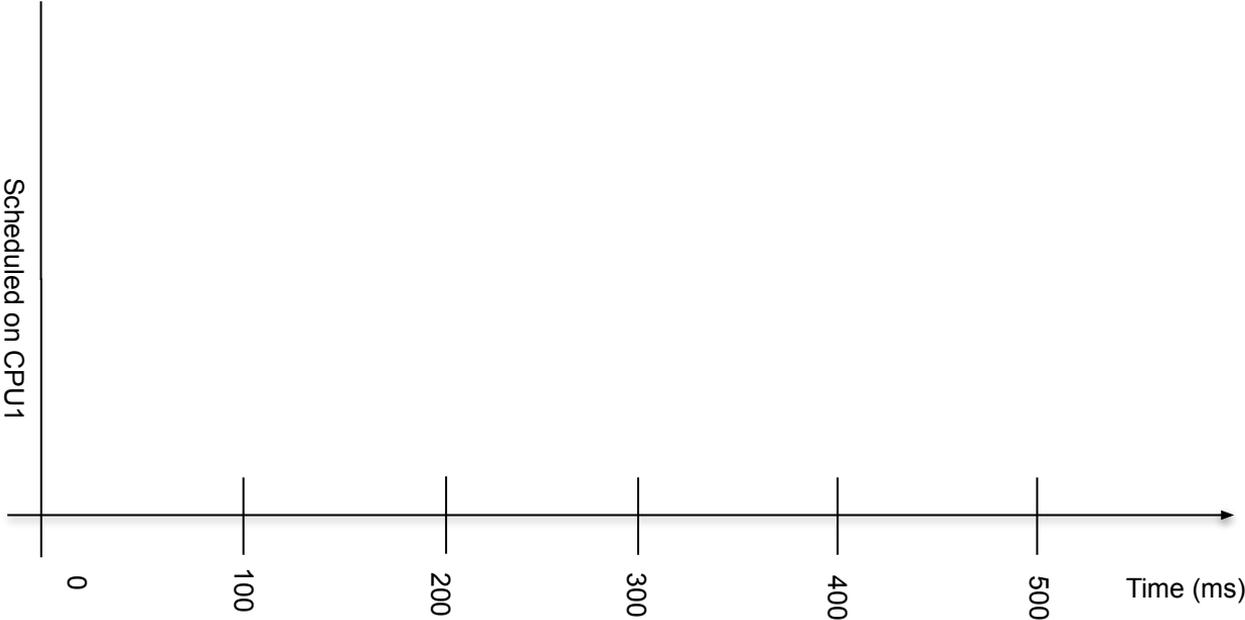
nr\_throttled is number of those periods the application was throttled

Throttled percentage

$$\text{throttled}\% = \frac{\Delta nr\_throttled}{\Delta nr\_periods}$$

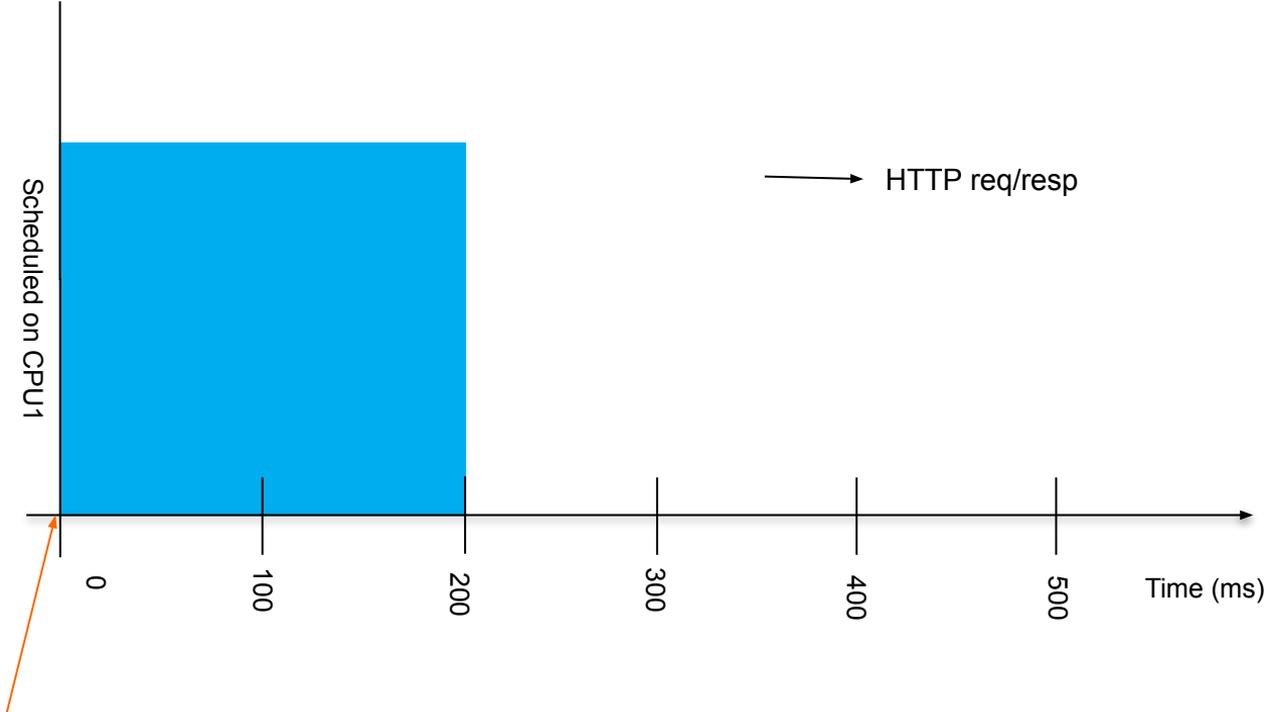
# Conceptual Model: Unconstrained

CPU time required 200ms



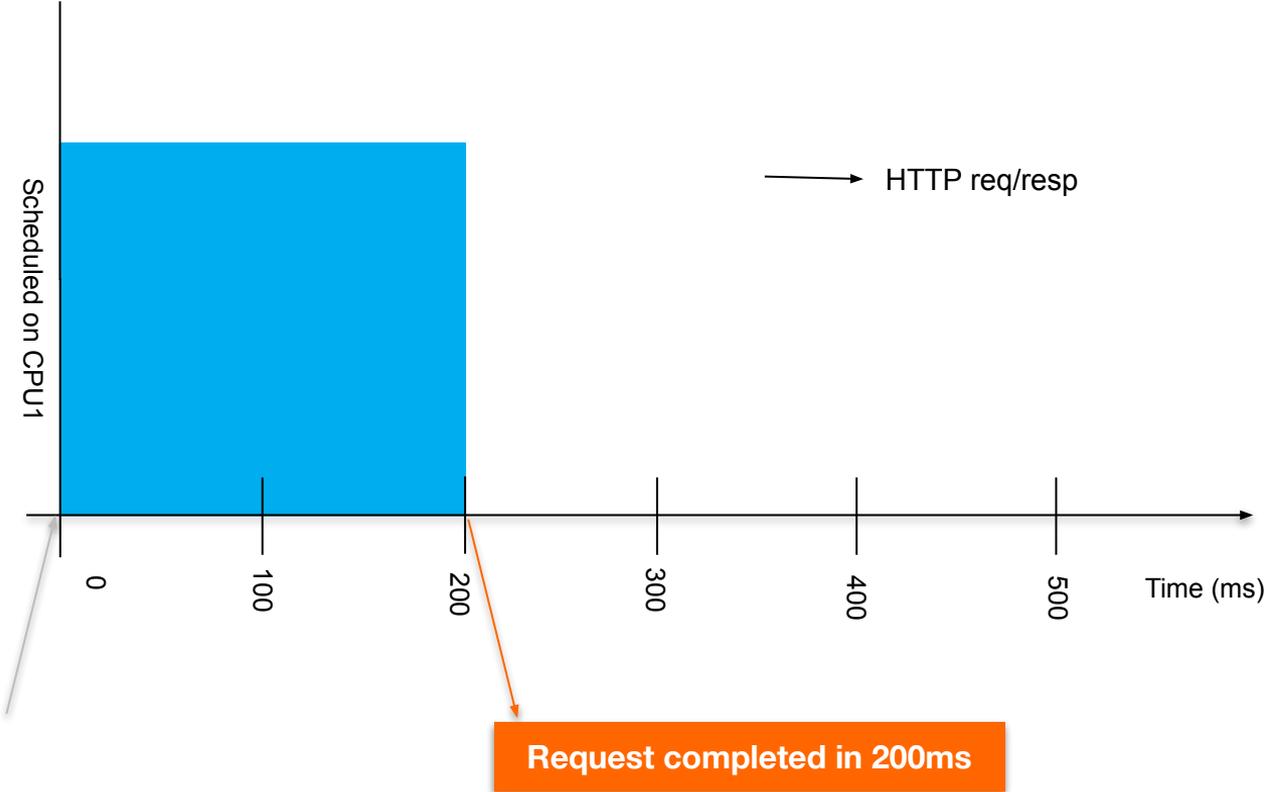
# Conceptual Model: Unconstrained

CPU time required 200ms



# Conceptual Model: Unconstrained

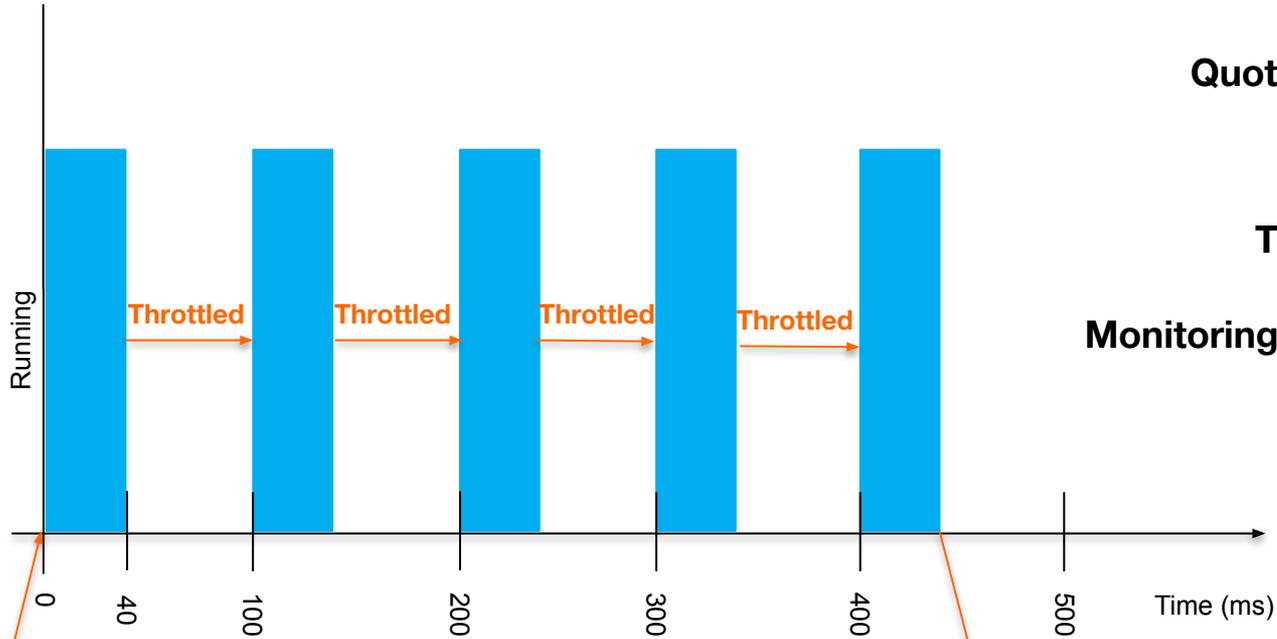
CPU time required 200ms



# Conceptual Model: CFS Bandwidth Control

CPU time required 200ms  
CPU quota .4 CPU  
Quota per period 40ms

CPU Usage 200ms  
Throttle Time 240ms  
Throttled % 80%  
Monitoring CPU Usage .4 CPU



Request

Request completed in 440ms

# Roadmap



CPU Limit Basics



**The Problem**

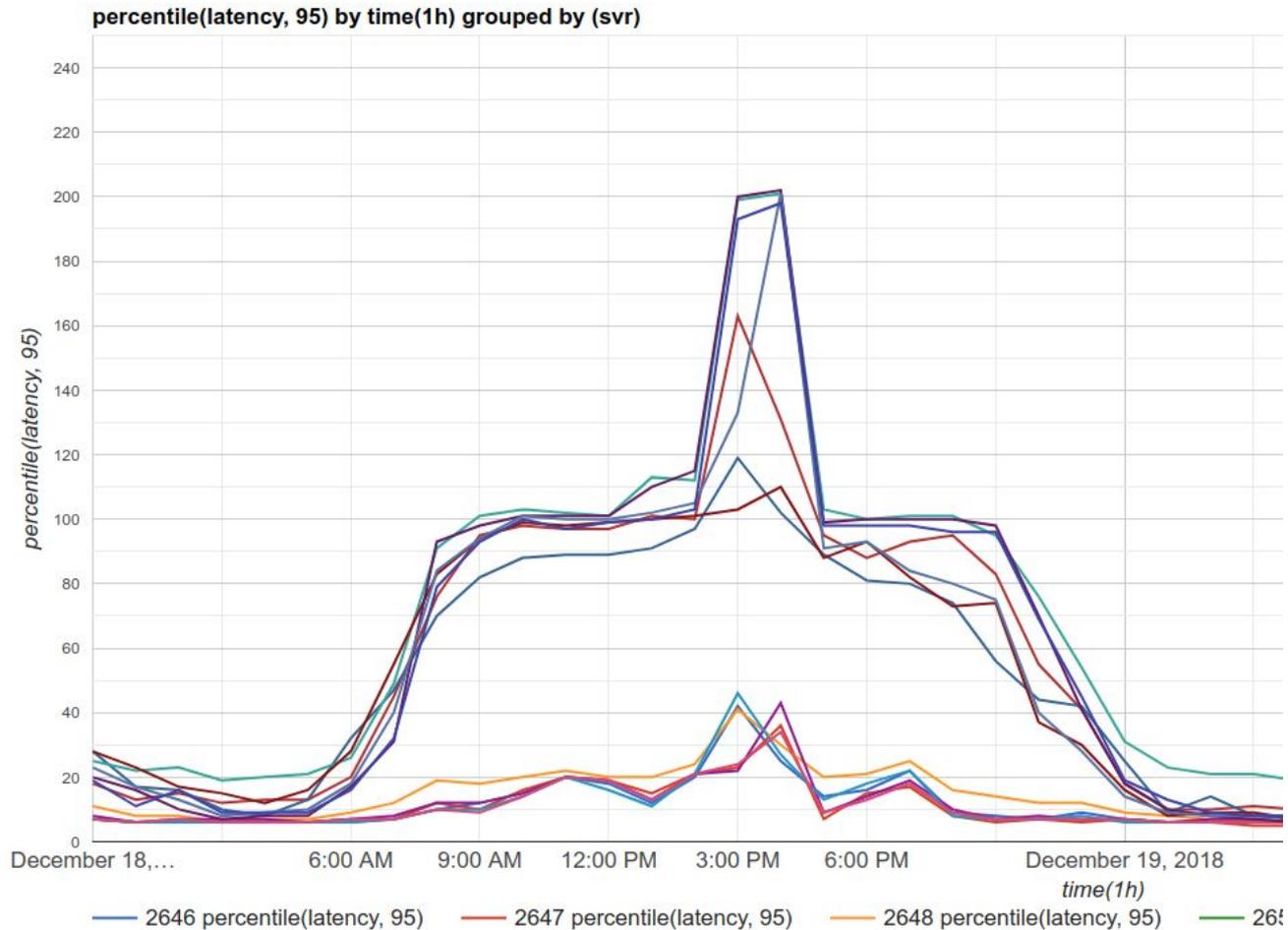


Reproducing the  
Problem



Solution and  
Workarounds

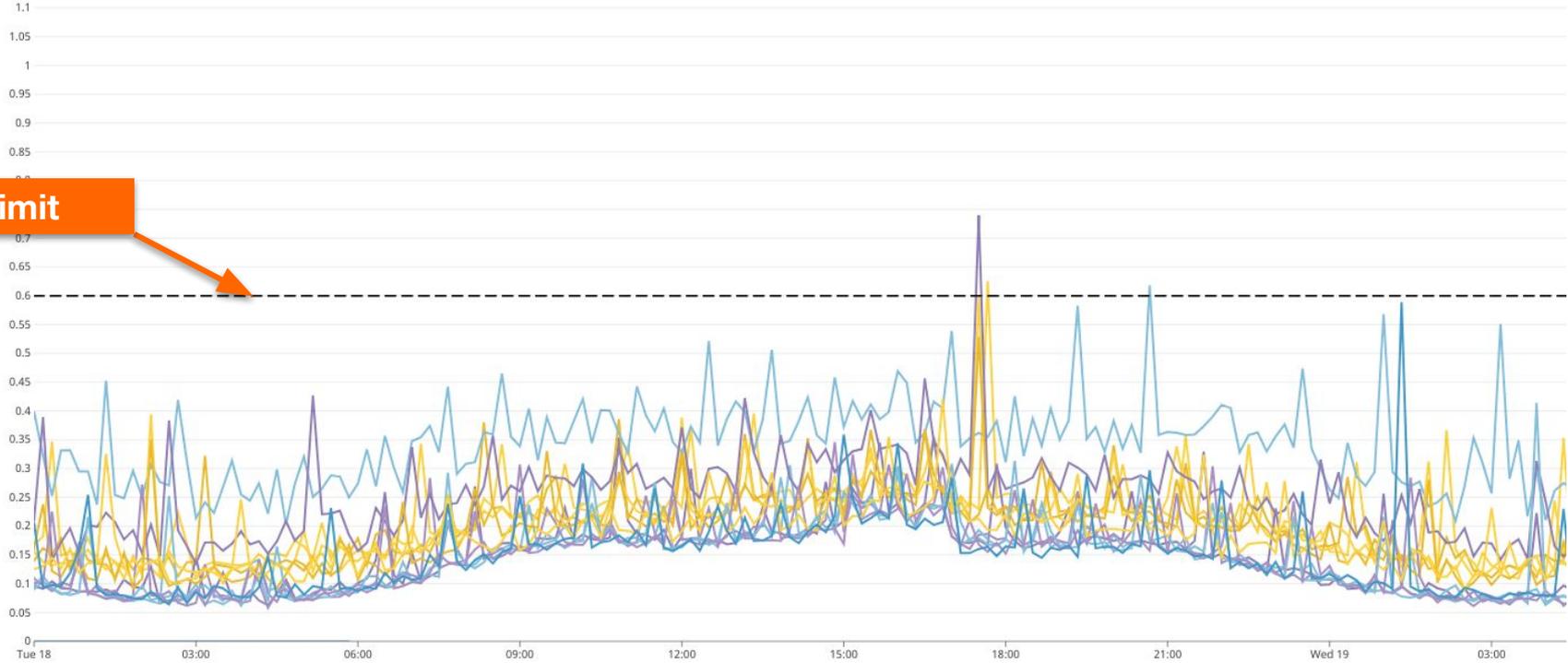
# Latency Issues



# CPU Usage

CPU Cores Used

47 h Dec 18, 12:00 am - Dec 19, 11:28 pm



Limit

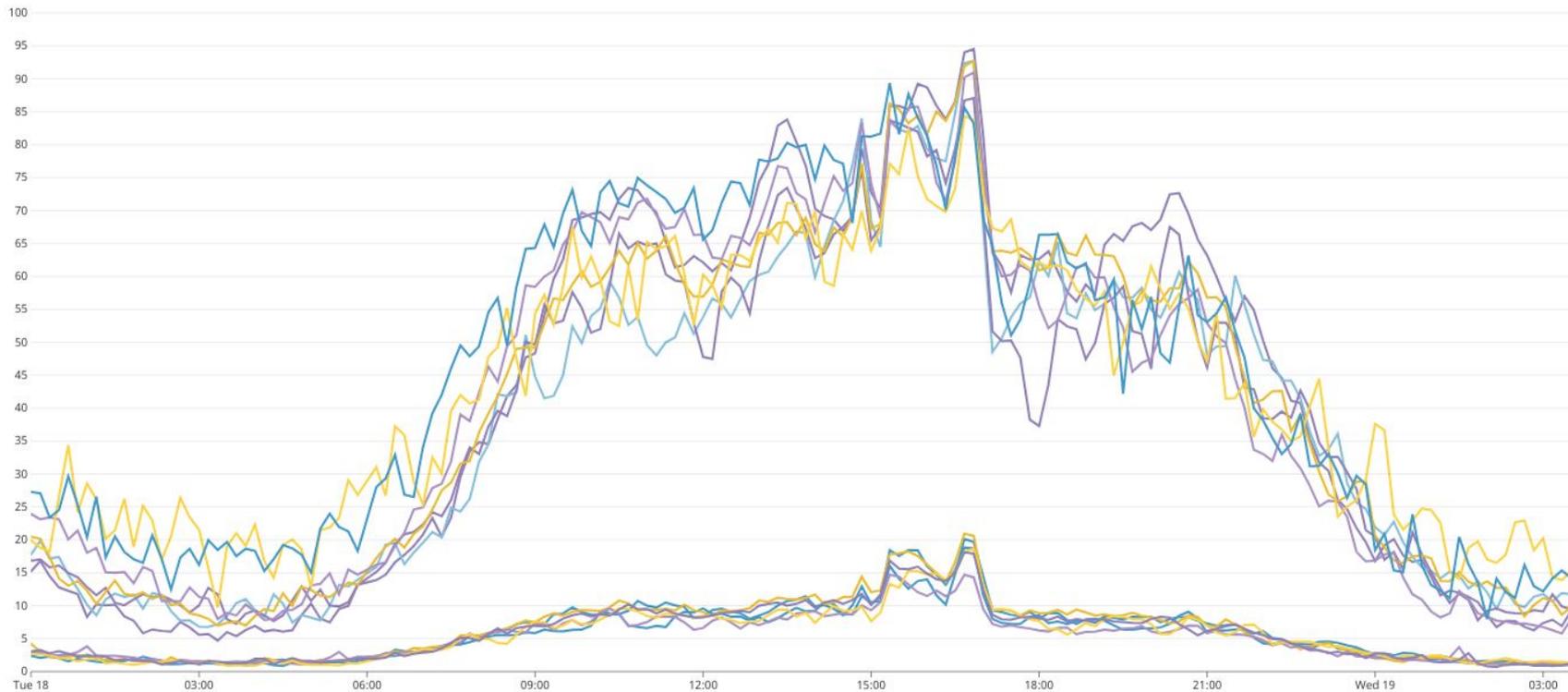


Filter series

# Throttling

ADVANCED: Percentage full quota usage by 100ms buckets (higher is worse)

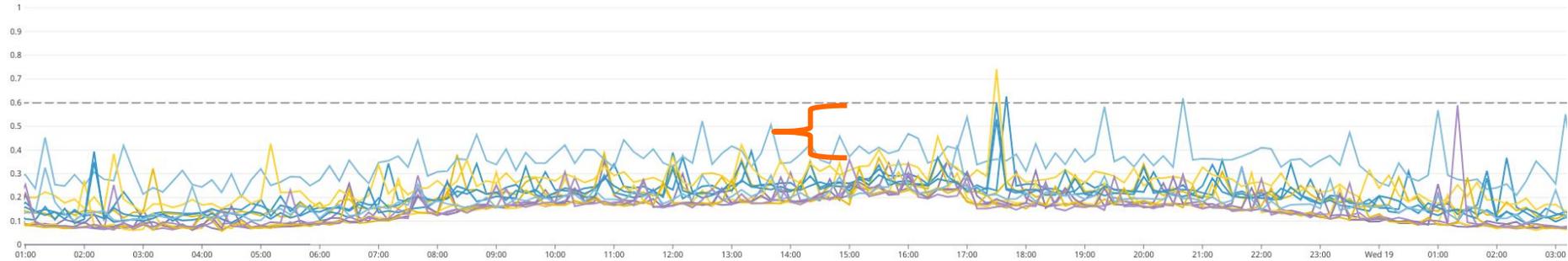
47 h Dec 18, 12:00 am - Dec 19, 11:28 pm



# Low CPU Usage but High Throttling?

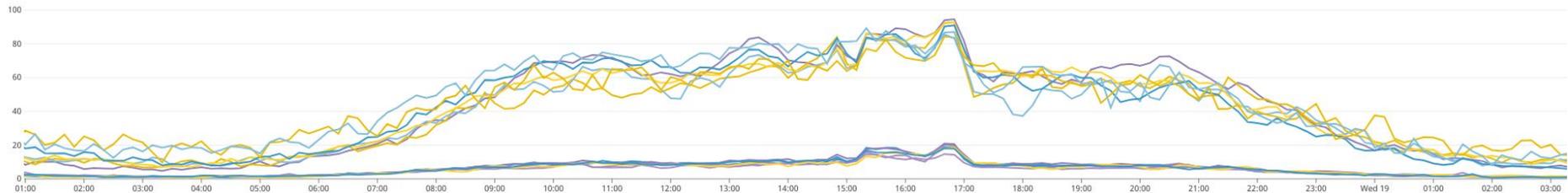
CPU Cores Used

29 h Dec 18, 12:58 am - Dec 19, 5:31 am



ADVANCED: Percentage full quota usage by 100ms buckets (higher is worse)

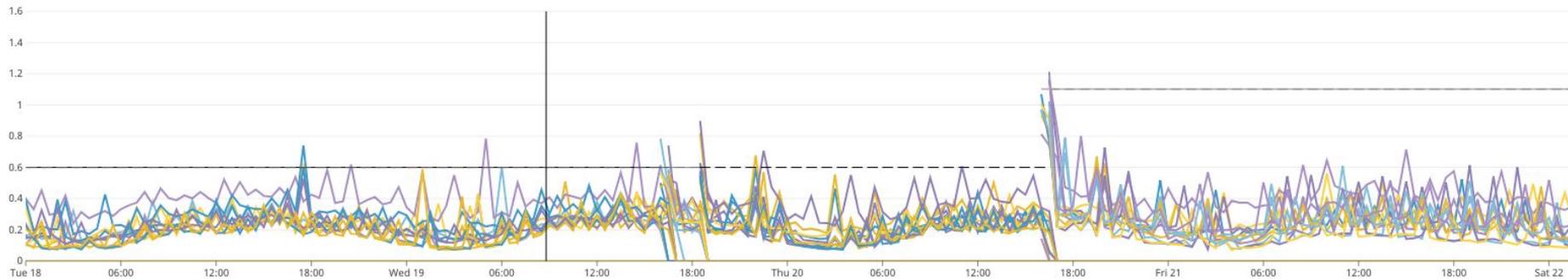
29 h Dec 18, 12:58 am - Dec 19, 5:31 am



# After increasing Limit

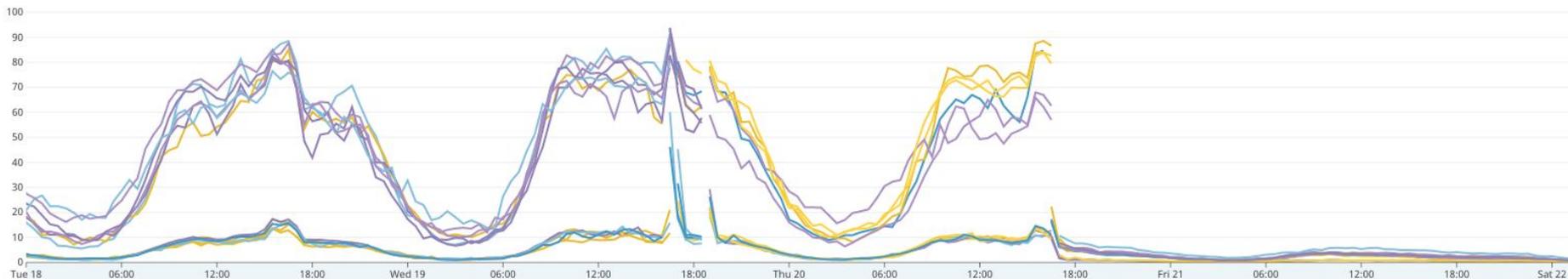
CPU Cores Used

5 d Dec 18, 12:00 am - Dec 22, 11:59 pm

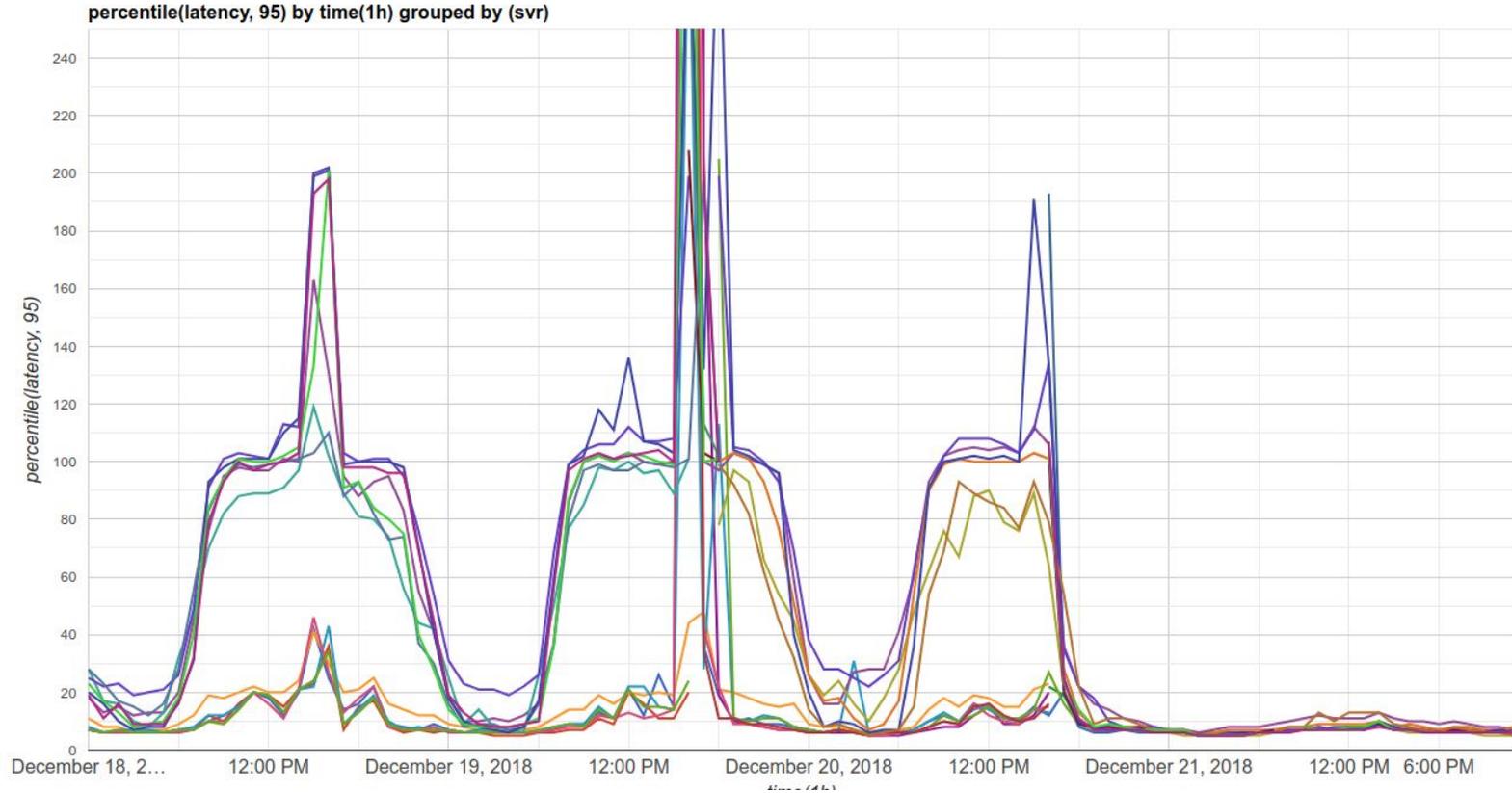


ADVANCED: Percentage full quota usage by 100ms buckets (higher is worse)

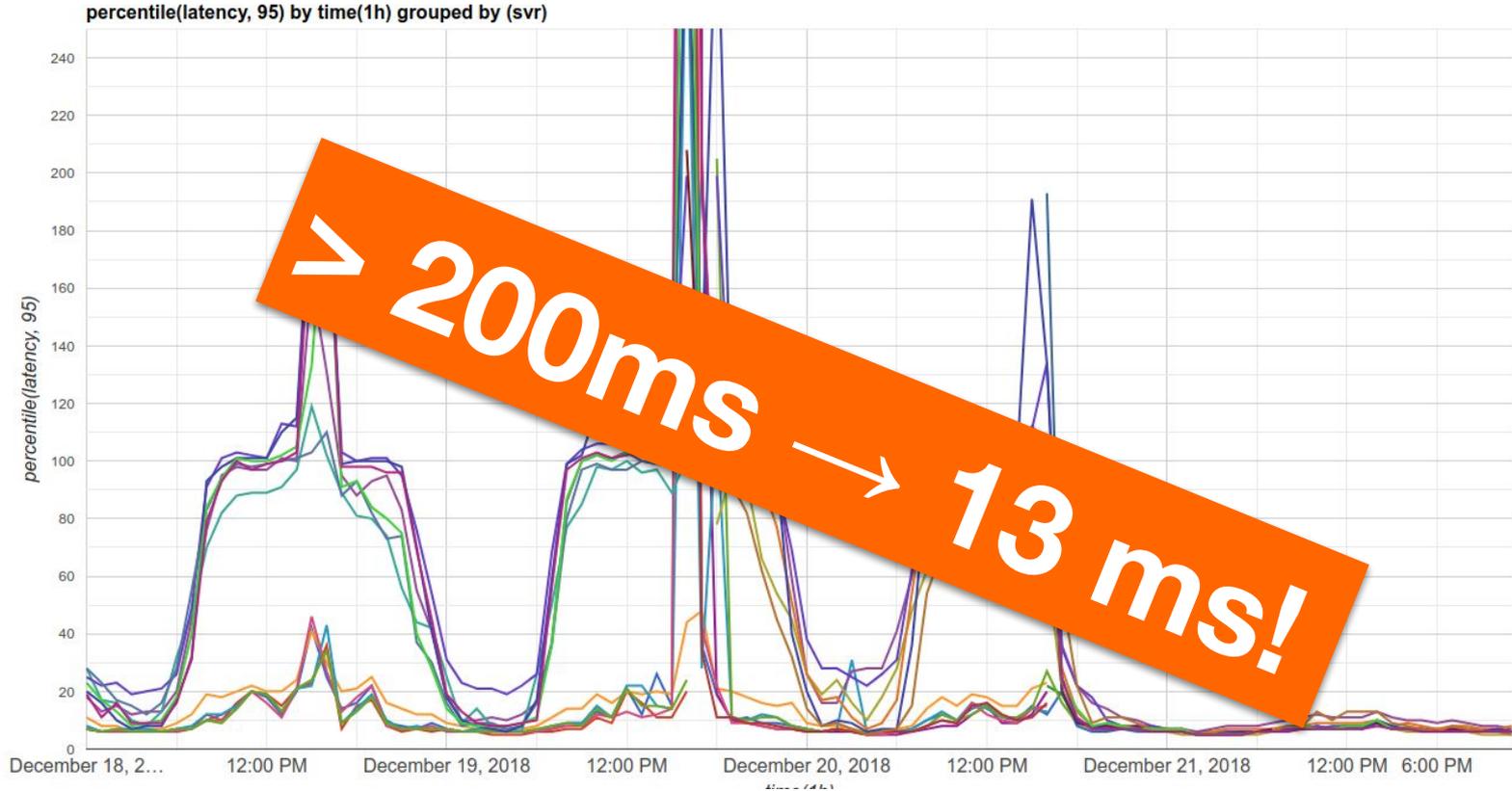
5 d Dec 18, 12:00 am - Dec 22, 11:59 pm



# After increasing Limit



# After increasing Limit



# What we know

# What we know

## → Workarounds

- + Increasing CPU quota mitigates throttling

# What we know

## → Workarounds

- + Increasing CPU limit mitigates throttling

## → Possible root causes

- + High Core Count ?
- + CPU architecture ?
- + Kernel version ?
- + Spectre-meltdown mitigations ?

# Roadmap



CPU Limit Basics



The Problem



**Reproducing  
the Problem**



Solution and  
Workarounds

## Reproducing

+ **ab**

+ **stress-ng**

+ <bash>

```
for (( i=1 ; i <= 1000 ; i++ )) ; do
```

```
    curl -s http://127.0.0.1:8888/info/healthcheck 2>&1 >/dev/null &
```

```
    sleep .005s
```

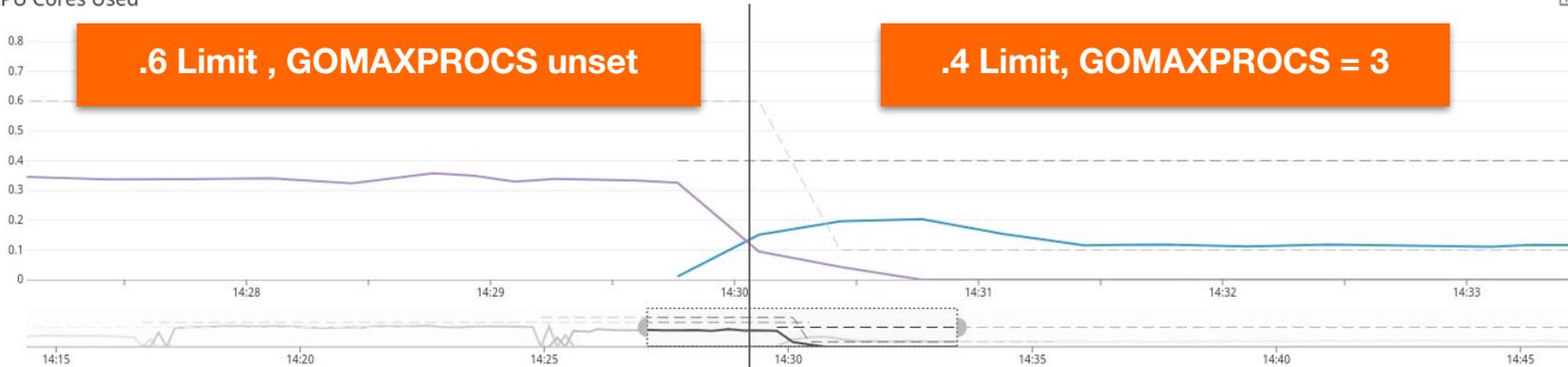
```
done
```

```
# Report amount of throttling
```

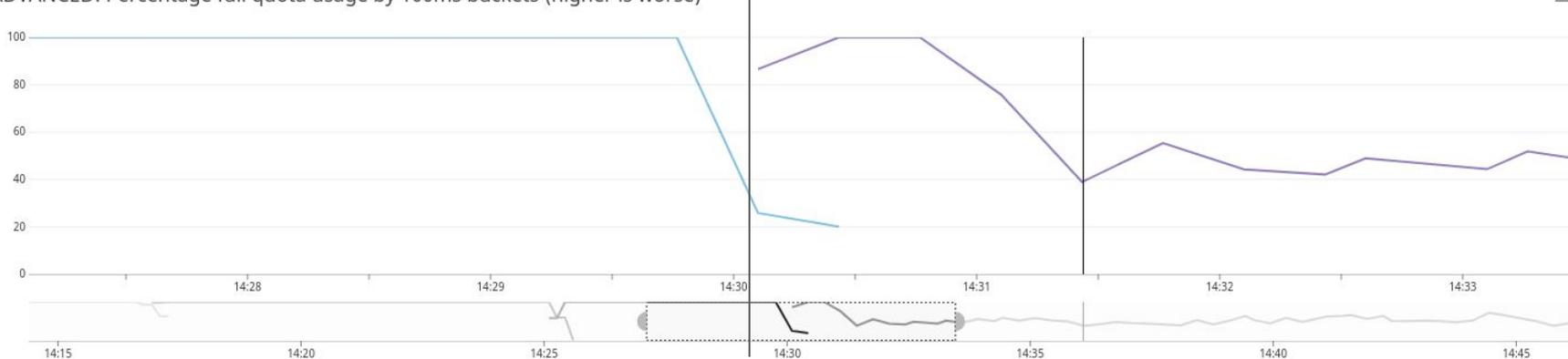
```
</bash>
```

# GOMAXPROCS

CPU Cores Used



ADVANCED: Percentage full quota usage by 100ms buckets (higher is worse)



## What We Know

### Workarounds

- + Increase CPU quota
- + Decrease number of threads in the application
  - + Golang - set GOMAXPROCS
  - + Java - move to newer JVMs that are cgroup aware.
- + Move from fractional to whole cpu shares.

### To Do

- + Create Custom Reproducer
- + Fix Kernel Scheduler

The background is a vibrant orange color with a complex, abstract pattern. It features several sets of concentric circles and arcs in various shades of orange, creating a sense of depth and movement. There are also some rectangular and semi-circular shapes, some with a fine dotted texture, scattered across the composition. The overall effect is modern and dynamic.

# The Reproducer: Fibtest

# Fibtest

## Multithreaded calculation of the Fibonacci sequence

- + Fast threads – calculate as fast as possible
- + Slow threads – calculate 100 iterations, then sleep for 10ms
- + Each thread is pinned to it's own CPU
- + <https://github.com/indeedeng/fibtest>

# Running Fibtest

```
[fibtest]$ ./runfibtest 1  
Iterations Completed(M): 1452  
Throttled for: 10  
CPU Usage (msecs) = 539  
[fibtest]$ ./runfibtest 16  
Iterations Completed(M): 1380  
Throttled for: 11  
CPU Usage (msecs) = 530  
[fibtest]$ ./runfibtest 88  
Iterations Completed(M): 275  
Throttled for: 11  
CPU Usage (msecs) = 183
```

# Running Fibtest

```
[fibtest]$ ./runfibtest 1
Iterations Completed(M): 1452
Throttled for: 10
CPU Usage (msecs) = 539
[fibtest]$ ./runfibtest 16
Iterations Completed(M): 1380
Throttled for: 11
CPU Usage (msecs) = 530
[fibtest]$ ./runfibtest 88
Iterations Completed(M): 275
Throttled for: 11
CPU Usage (msecs) = 183
```



# Running Fibtest

```
[fibtest]$ ./runfibtest 1
Iterations Completed(M): 1452
Throttled for: 10
CPU Usage (msecs) = 539
[fibtest]$ ./runfibtest 16
Iterations Completed(M): 1380
Throttled for: 11
CPU Usage (msecs) = 530
[fibtest]$ ./runfibtest 88
Iterations Completed(M): 275
Throttled for: 11
CPU Usage (msecs) = 183
```



## Running Fibtest

```
[fibtest]$ ./runfibtest 1
Iterations Completed(M): 1452
Throttled for: 10
CPU Usage (msecs) = 539
[fibtest]$ ./runfibtest 16
Iterations Completed(M): 1380
Throttled for: 11
CPU Usage (msecs) = 530
[fibtest]$ ./runfibtest 88
Iterations Completed(M): 275
Throttled for: 11
CPU Usage (msecs) = 183
```

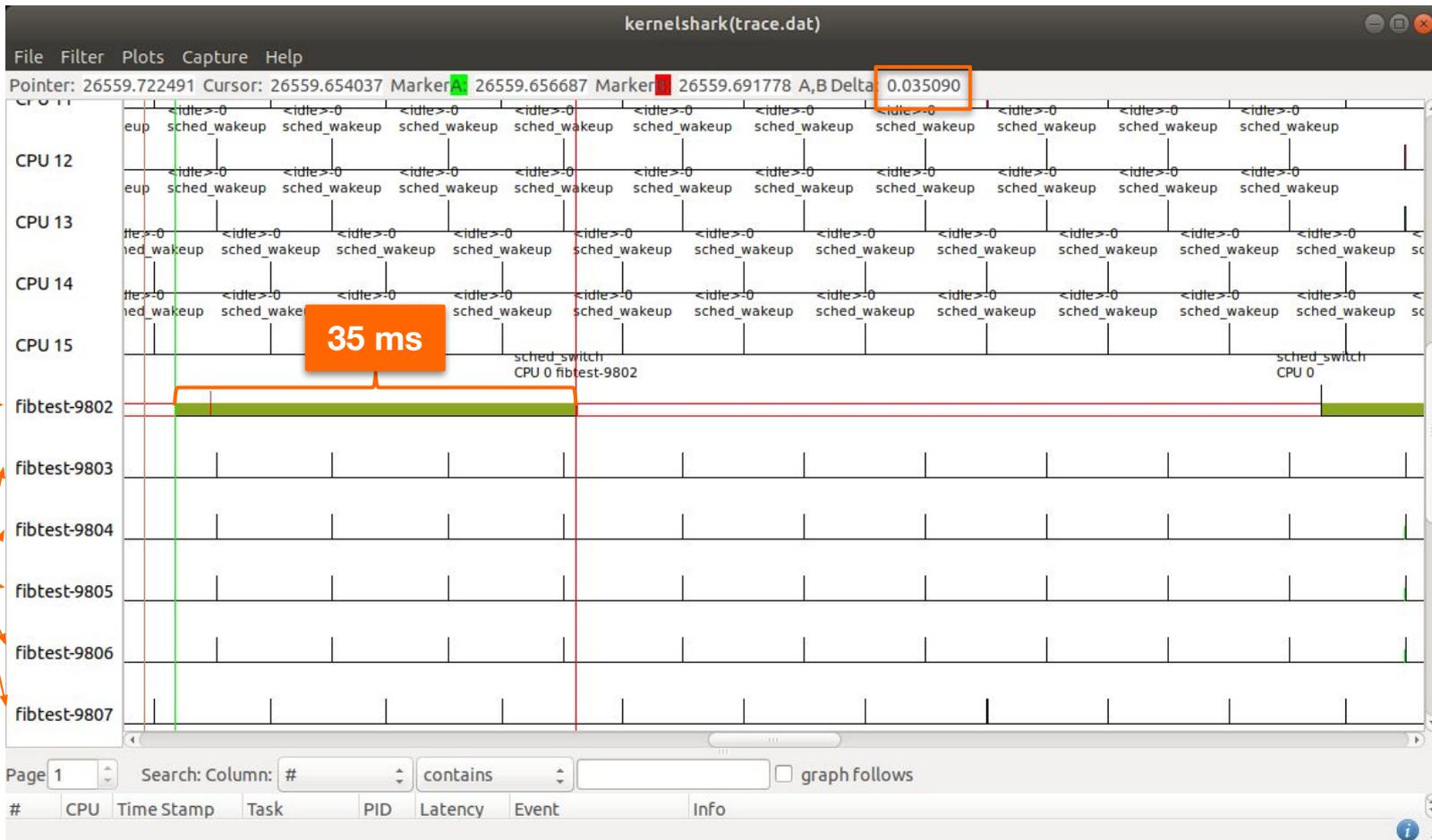
~3x

# Fibtest Tracing

```
$ sudo trace-cmd record -e 'sched_wakeup*' -e sched_switch -F -c ./runfibtest 16
```



# ftrace/kernelshark of Fibtest on 4.18-5.3.8



The background is a vibrant orange color with a complex, abstract pattern. It features several sets of concentric circles and arcs in various shades of orange, some solid and some with a dotted texture. The shapes are layered and overlapping, creating a sense of depth and movement. The overall aesthetic is modern and geometric.

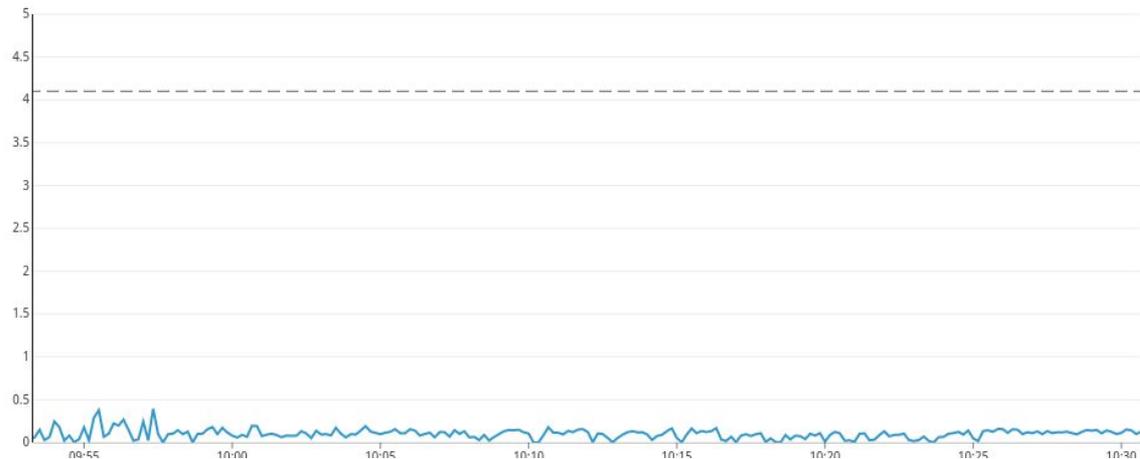
# The Causal Commit

512ac999

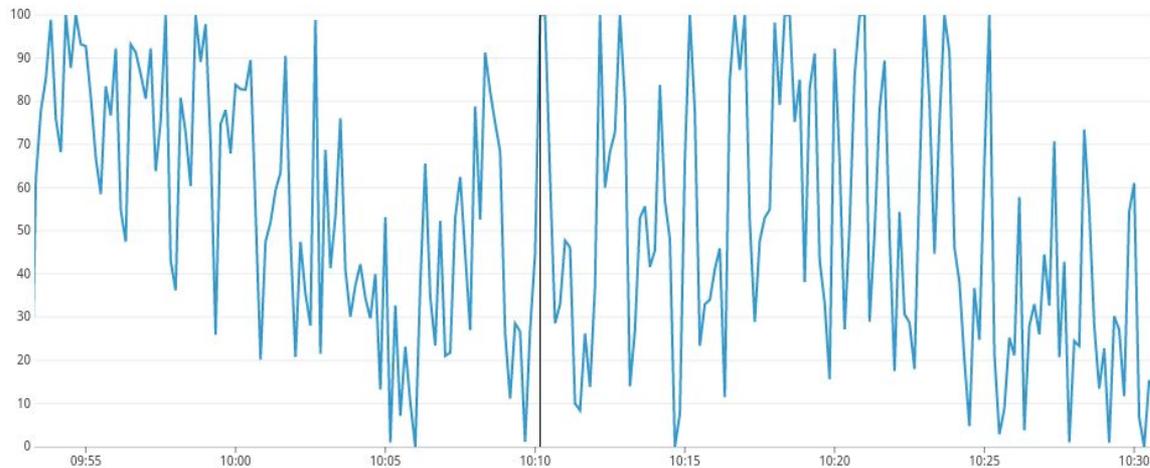
→ **Fix for inadvertent throttling due to clock-drift**

# 512ac999 Clock-drift problem

CPU Cores Used



Percentage of 100ms Periods Throttled.



## 512ac999

- Fix for inadvertent throttling due to clock-drift
- **Fixed per-cpu quota to expire on period boundaries**

# Real World vs. Conceptual Model

- **Multiple CPUs**
- **Many threads \*(sometimes thousands)**
- **Cores run at different speeds - Use performance mode**
- **Schedulers are hard**

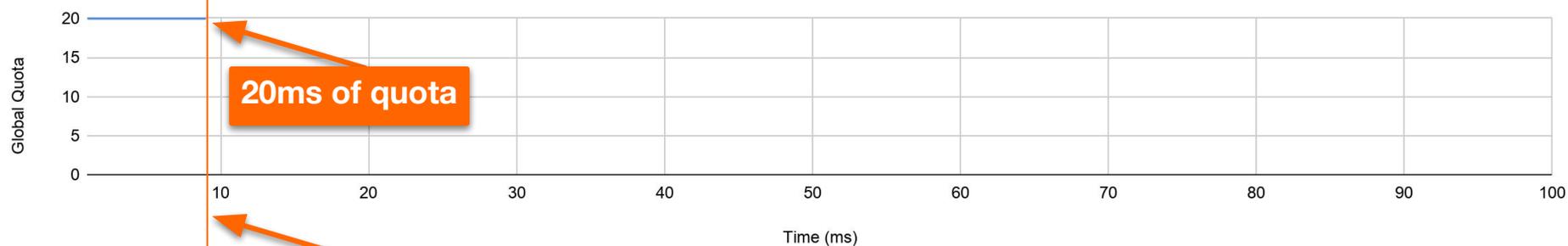
Quota is split into 5ms slices and assigned to individual CPUs

(1 CPU of quota = 100ms/period = 20 slices/period) = not enough for large machines

Per-cpu quota will expire if not used within a period

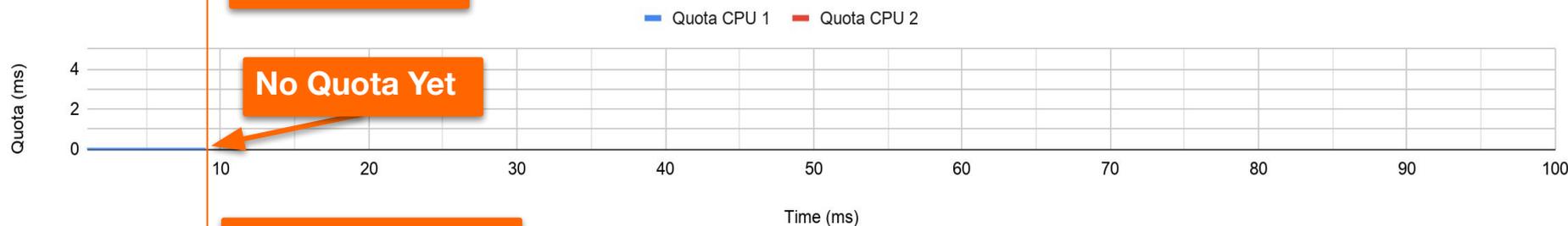
# Real World Example

Available Global Quota



20ms of quota

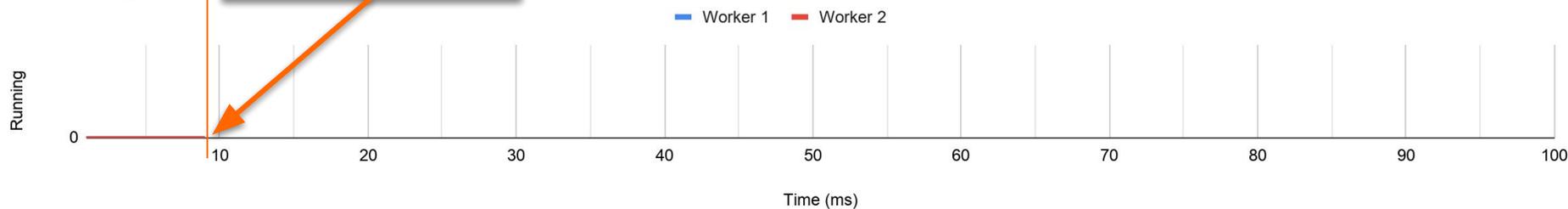
Per-CPU Quota



Current Time

No Quota Yet

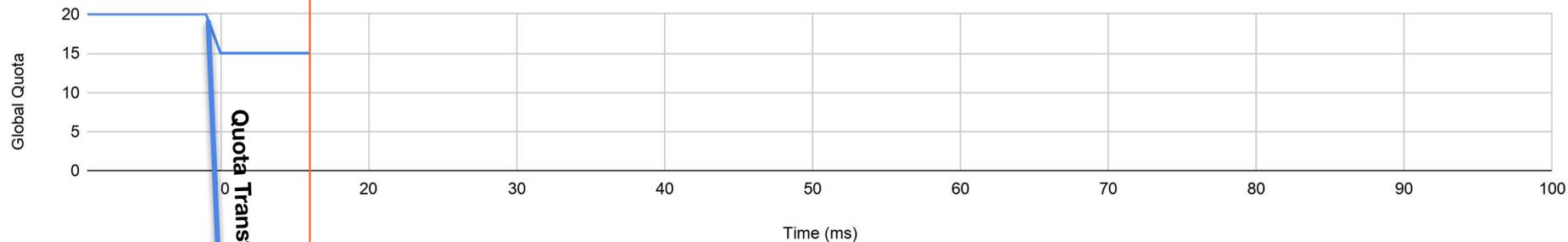
CPU Usage



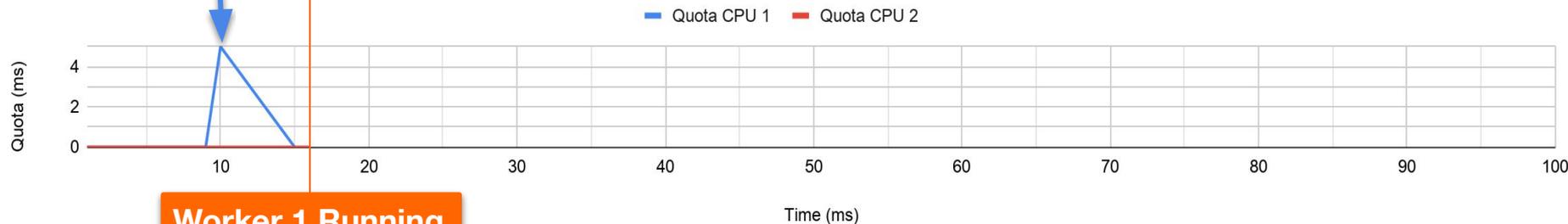
Nothing Running

# Real World Example

Available Global Quota

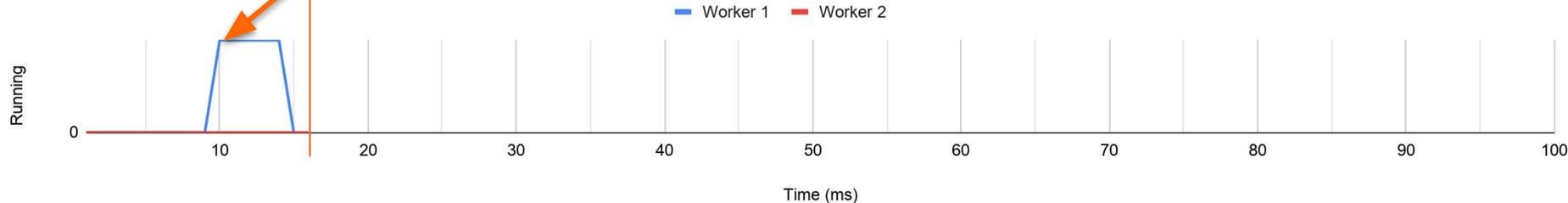


Per-CPU Quota



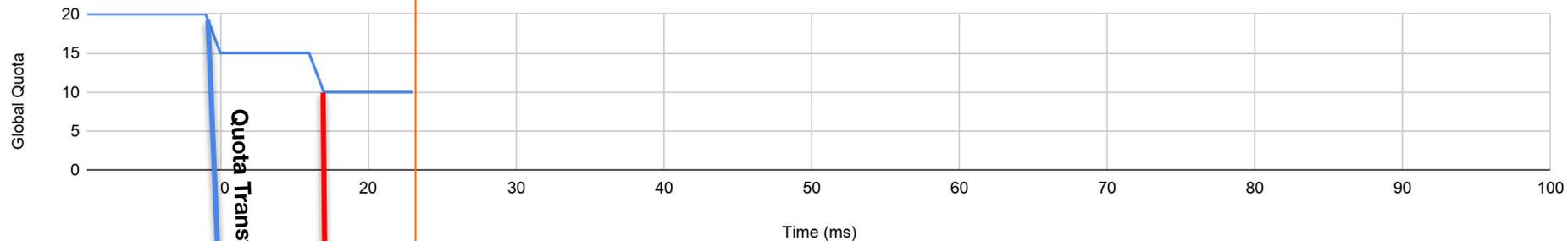
**Worker 1 Running**

CPU Usage

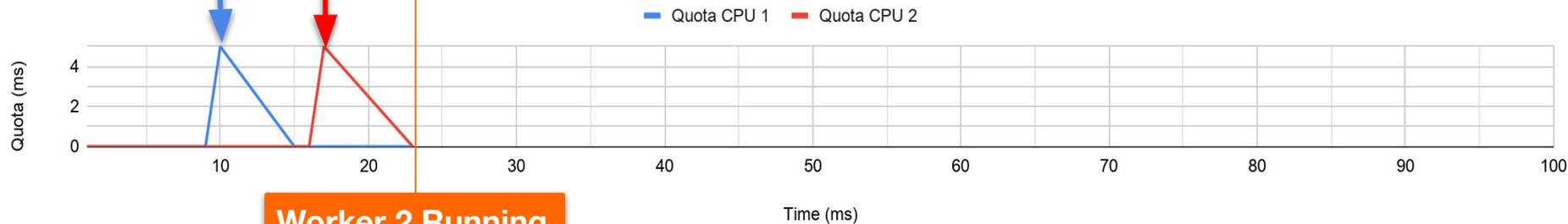


# Real World Example

Available Global Quota

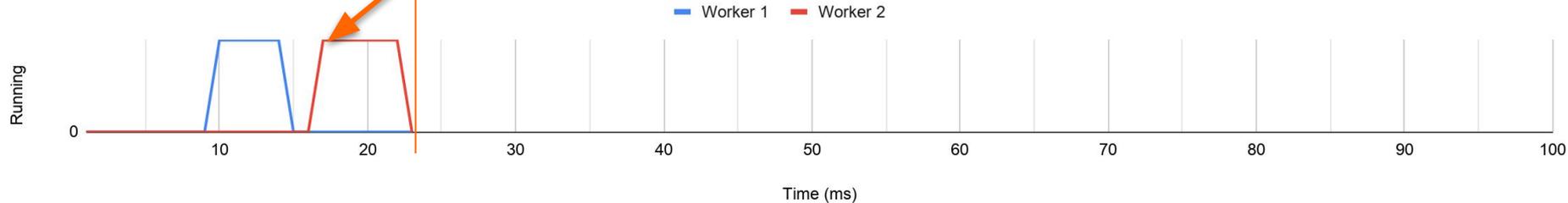


Per-CPU Quota



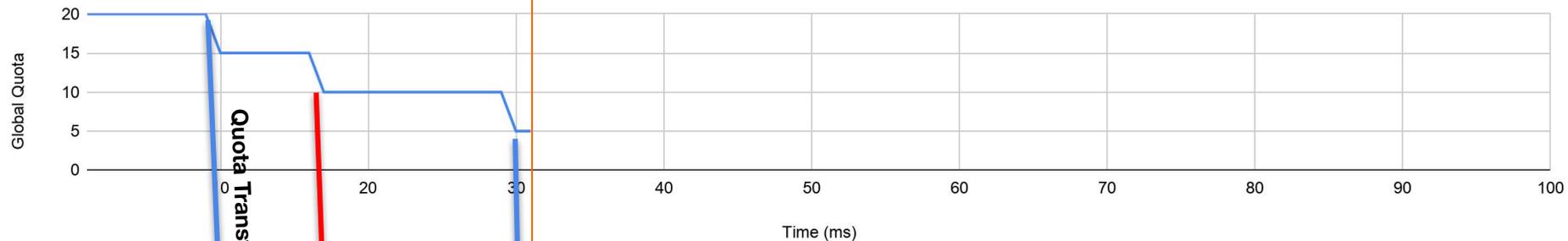
**Worker 2 Running**

CPU Usage

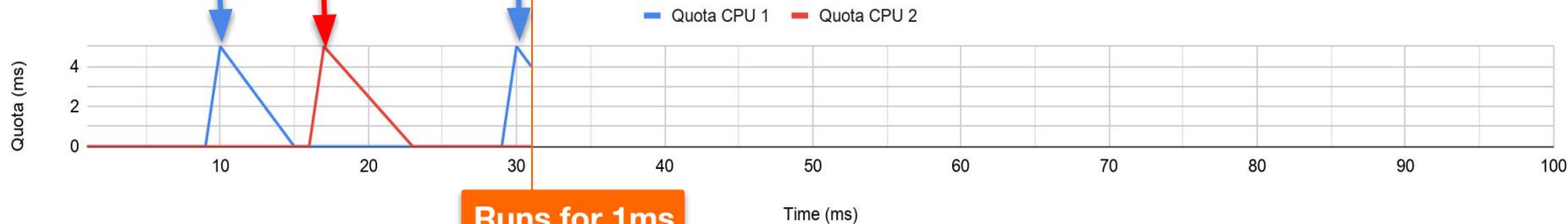


# Real World Example

## Available Global Quota



## Per-CPU Quota

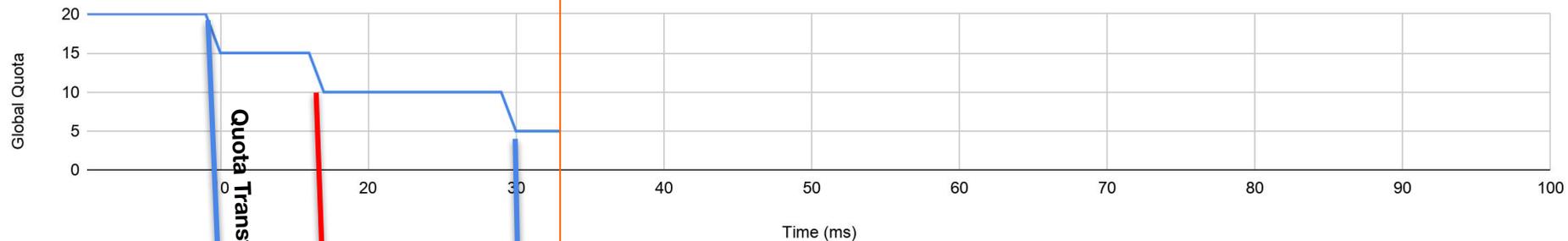


## CPU Usage

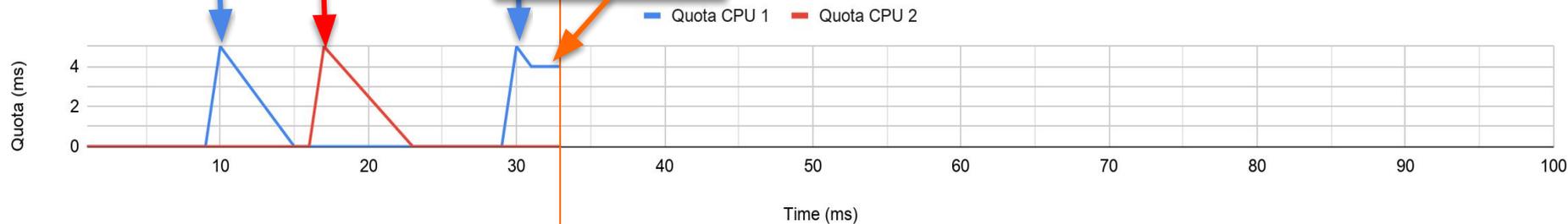


# Real World Example

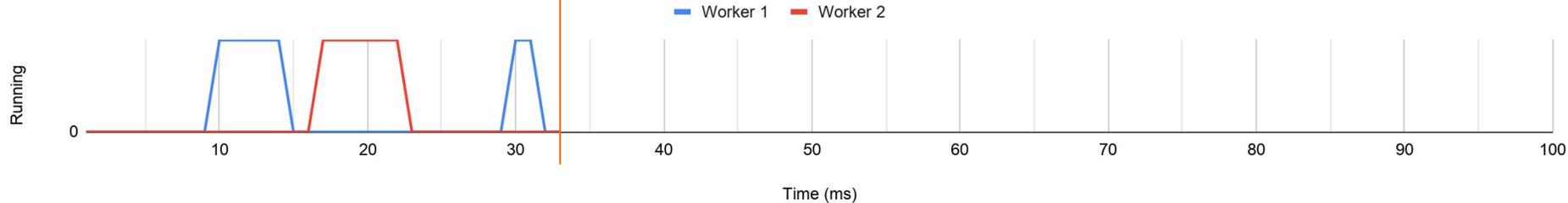
## Available Global Quota



## Per-CPU Quota

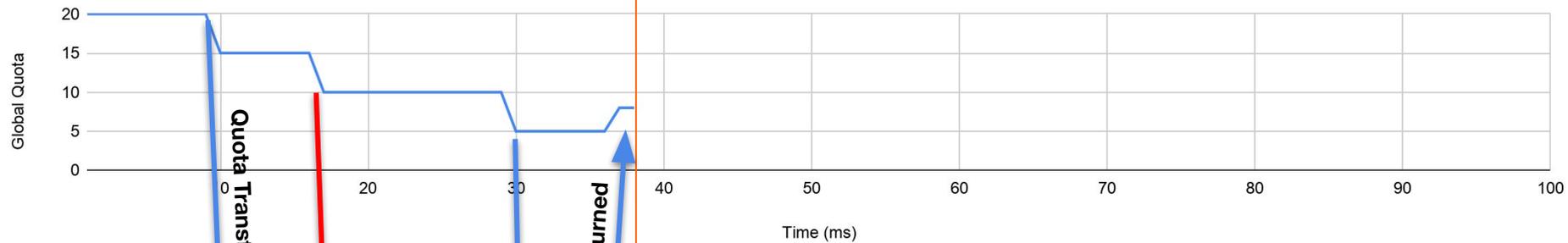


## CPU Usage

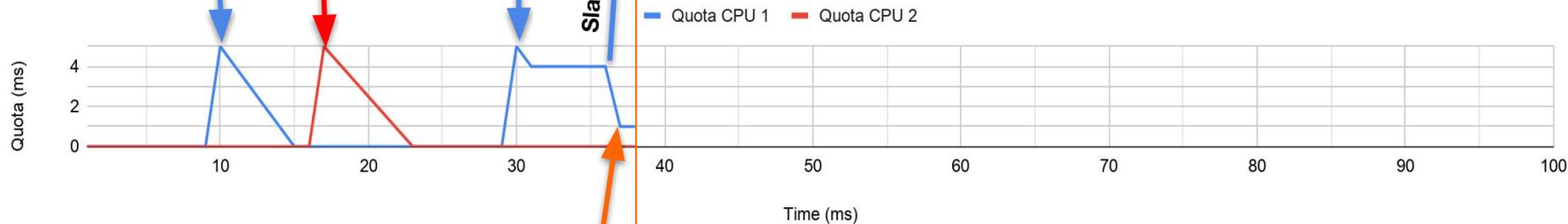


# Real World Example

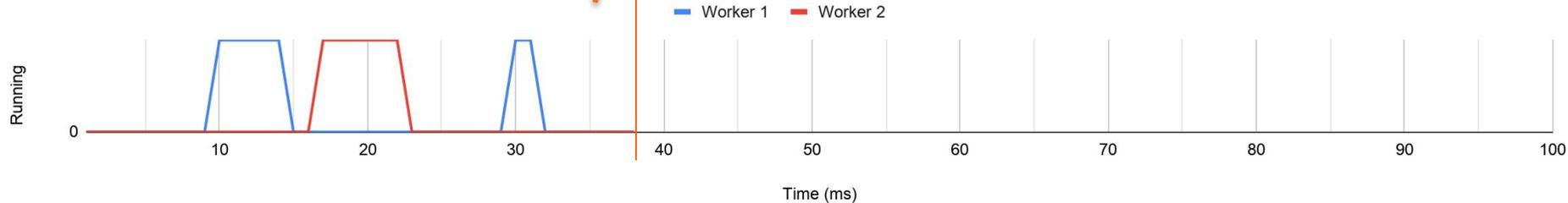
## Available Global Quota



## Per-CPU Quota



## CPU Usage

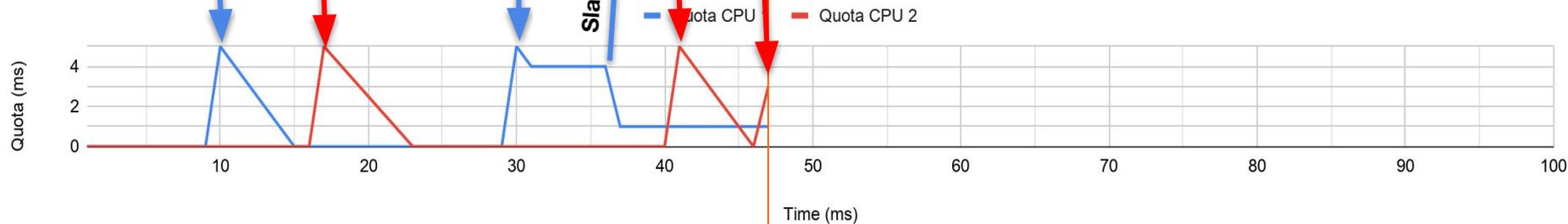


# Real World Example

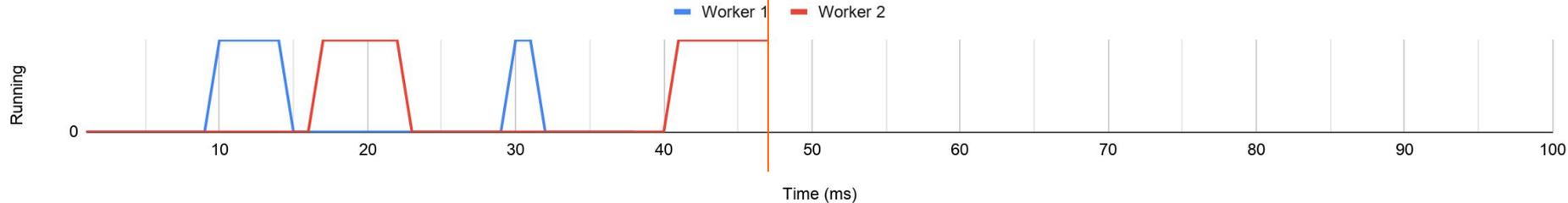
## Available Global Quota



## Per-CPU Quota

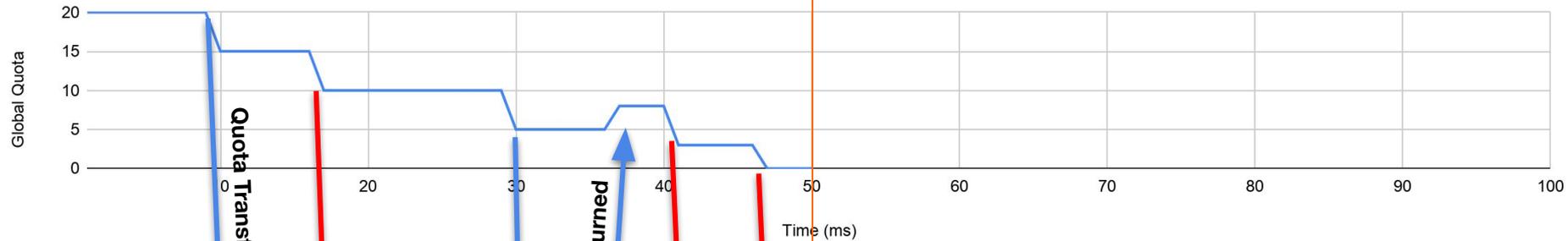


## CPU Usage

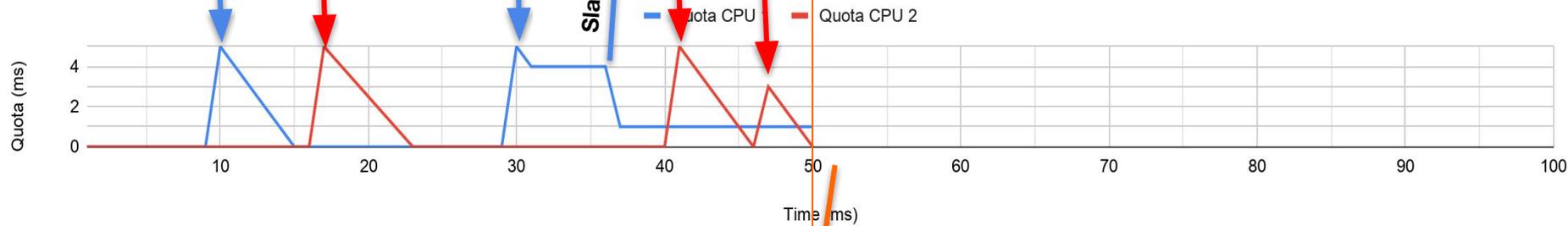


# Real World Example

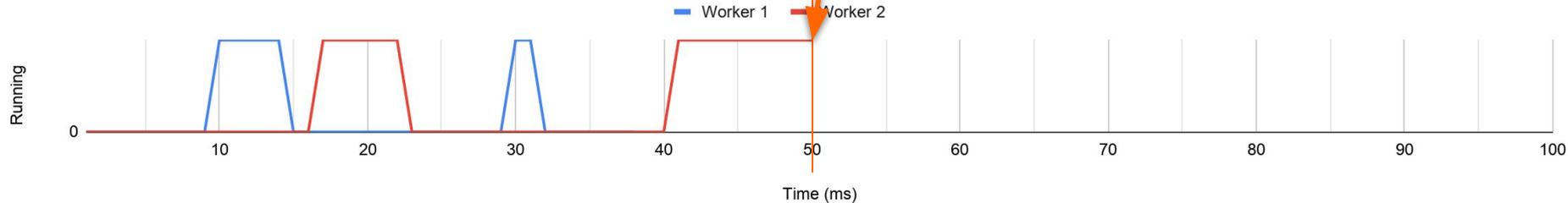
## Available Global Quota



## Per-CPU Quota

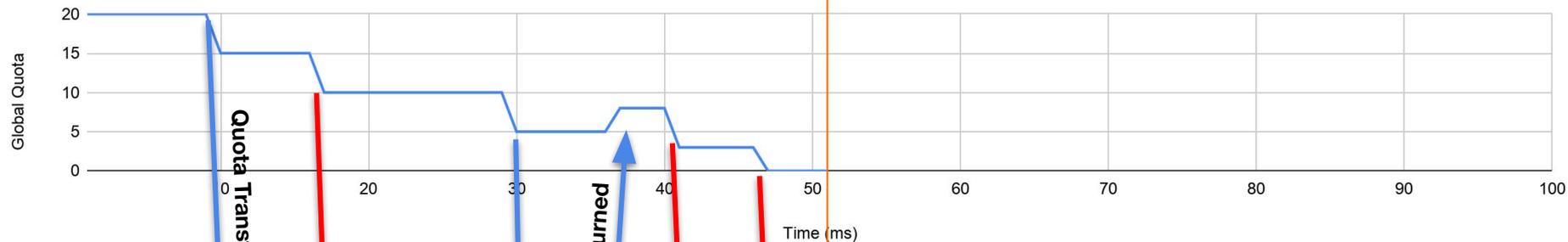


## CPU Usage



# Real World Example

Available Global Quota



Per-CPU Quota

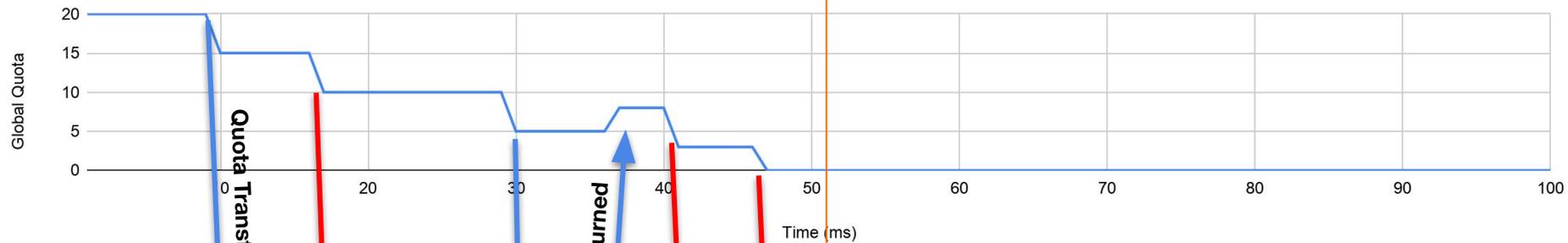


CPU Usage

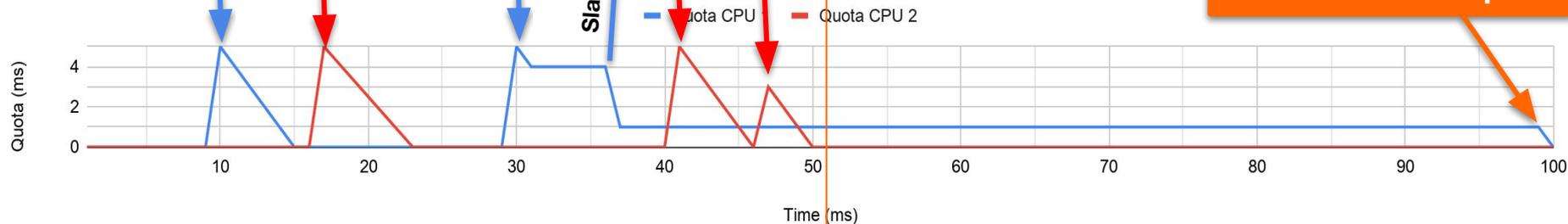


# Real World Example

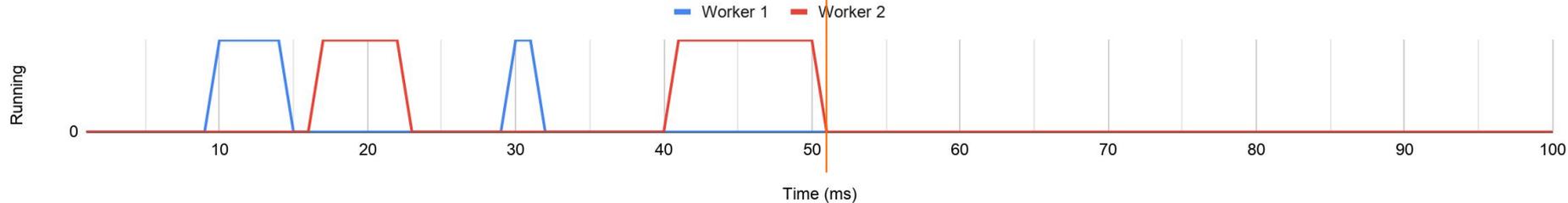
## Available Global Quota



## Per-CPU Quota



## CPU Usage



# IMPACT

$$1\text{ms}/100\text{ms} * (88 \text{ CPUs}-1) = 87\text{ms}/100\text{ms}$$

$$= 870\text{m} = .87 \text{ CPU}$$

# Roadmap



CPU Limit Basics



The Problem



Reproducing the  
Problem



**Solution and  
Workarounds**

## Possible Solutions

- **Remove 512ac999**
- **Burst bank / Rollover minutes**
- **Remove all per-CPU expiration logic**

## The Solution

**Remove all per-CPU expiration logic**

**→ 5 months of debate**

**→ 6 patch iterations**

## The Solution

**Commits: de53fd7aedb1 & 763a9ec06c40**

**→ Applied to 5.4 Kernel**

**→ linux-stable**

**◆ 4.14.154+, 4.19.84+, 5.3.9+**

**→ Distro kernels**

**◆ Ubuntu 5.3.0-24+**

**◆ Ubuntu 4.15.0-67+**

**◆ RHEL7 - kernel-3.10.0-1062.8.1.el7**

**◆ RHEL8.2 - WIP**

## The Solution

### Kernel 5.3.7

```
[fibtest]$ ./runfibtest 88  
Iterations Completed(M): 213  
Throttled for: 11  
CPU Usage (msecs) = 137
```



### Kernel 5.3.9

```
[fibtest]$ ./runfibtest 88  
Iterations Completed(M): 1316  
Throttled for: 10  
CPU Usage (msecs) = 482
```



**3x**

## Takeaways

- **Monitor your throttled %**
- **Upgrade your Kernels**
- **Use whole cpu quotas**
- **Increase quota where necessary**

**Questions?**



Dave Chiluk  
Linux Platform Software Engineer

@dchiluk

## Other Developments:

More developments:

Setting CFS Period ([GH #51135](#))

WIP:

Unset CFS quota with CPU sets ([GH #70585](#)) ([GH #75682](#))

## **Things I'd Like to see**

**Kernel: C-state aware quotas**

**Kernel: Burstable Cgroup CPU Limits**

**Kubernetes: Pod level Resource constraints**

**Kubernetes: Node Level CPU Overcommit**