# Improving Cassandra Client Load Balancing

Ammar Khaku
Joey Lynch

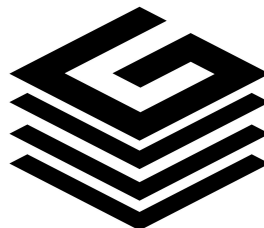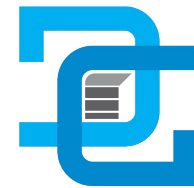**Speaker**



https://akhaku.com/

# Ammar Khaku

Senior Software Engineer
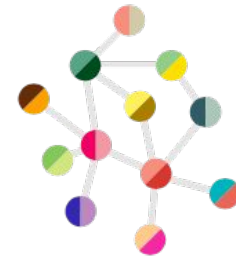Cloud Data Engineering at Netflix

Database clients, Java libraries

**Speaker**

# Joey Lynch

Senior Software Engineer
Cloud Data Engineering at Netflix
Cassandra Committer

Database shepherd and data wrangler

https://jolynch.github.io/

**Outline**

Load Balancing Background

Why Stateful Load Balancing is Special

Proposed Solution – Weighted Least Loaded

Experiments and Real World Results

# Goal: Upgrade to Datastax 4

Had some performance issues at scale with LoadBalancer and Throttler.

# (Un)Balance The Load

A quick crash course on [queueing theory](queueing theory) and load balancing

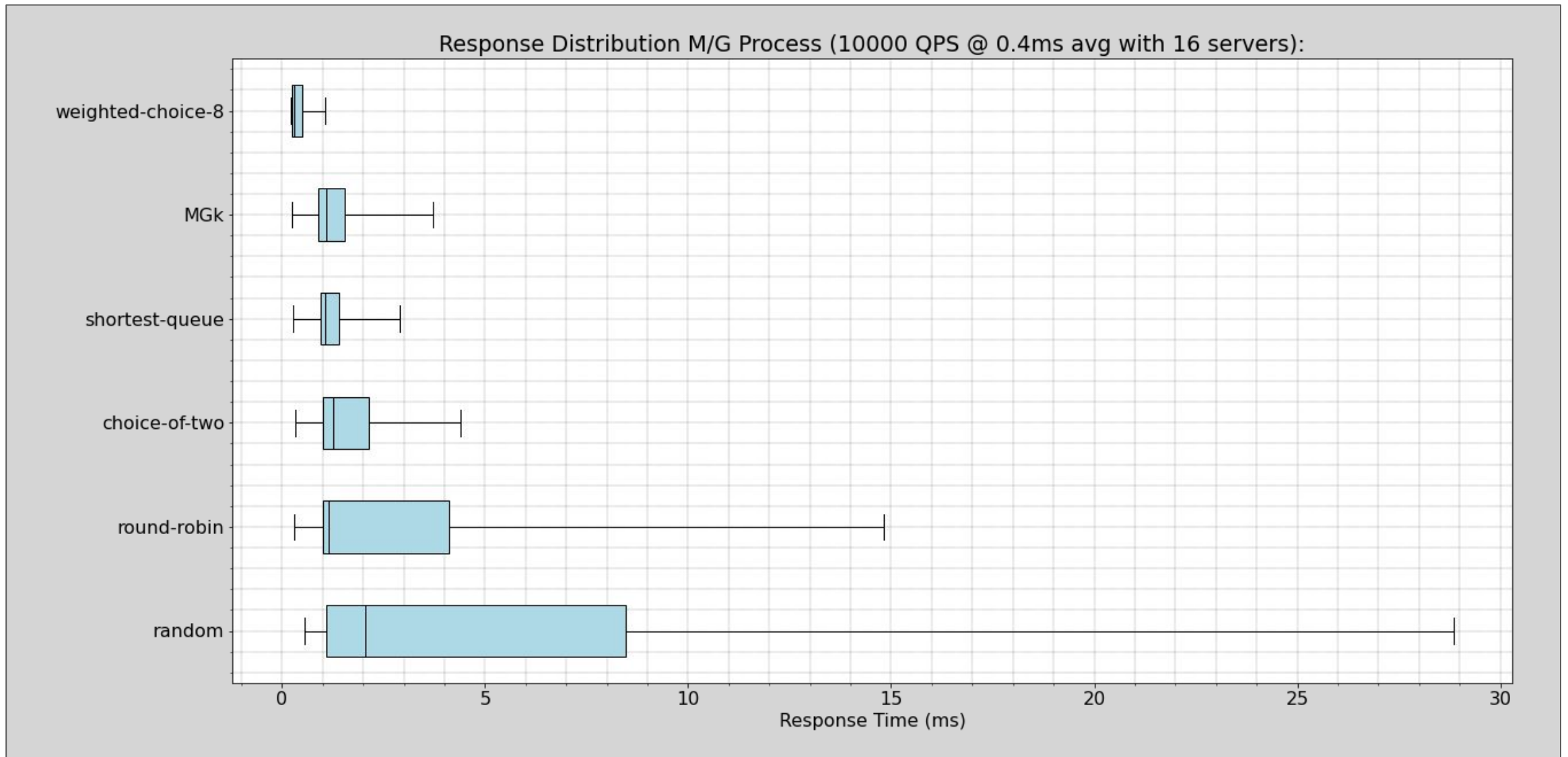**Best in class implementations**

# HAProxy, Nginx, Envoy

- Weighted **Round Robin**
- Weighted **Least Connection/Load**
- Weighted **Choice of N** (random/hash)

Netflix gRPC: **Random Choice of 2**

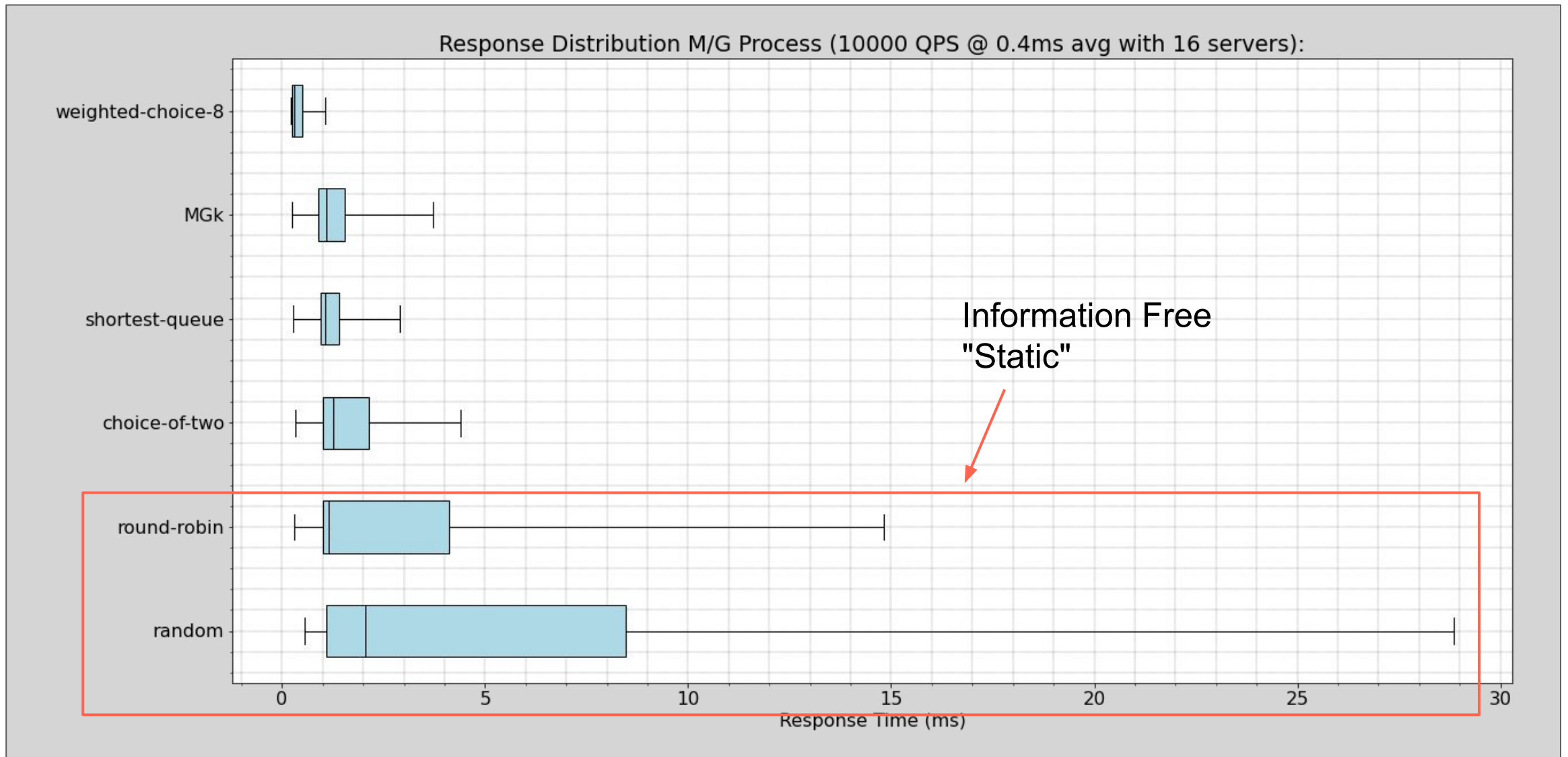Google uses **Random Subsetting** with weighted **Round Robin**
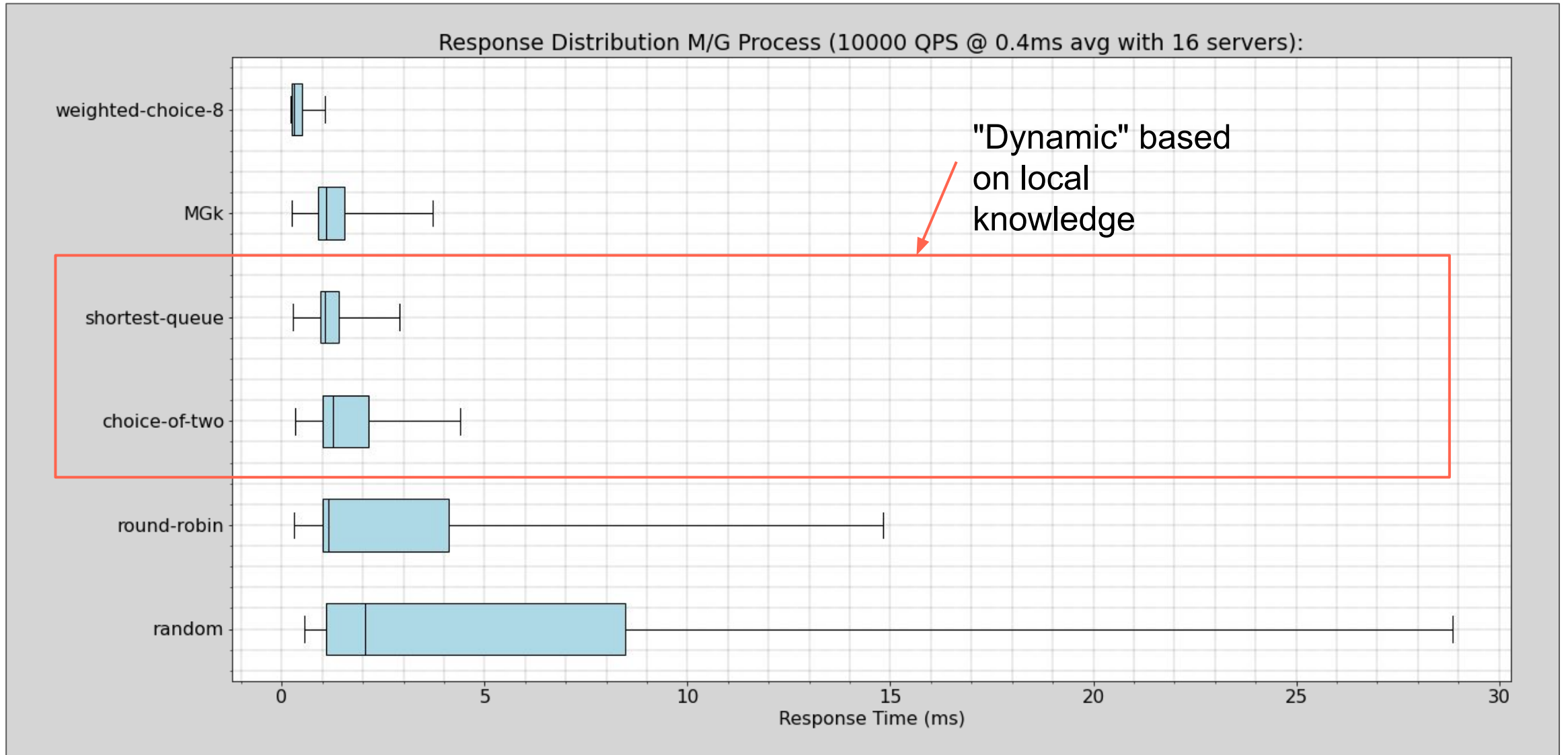
Many DB clients choose **Random**

# What to choose?



Response Distribution M/G Process (10000 QPS @ 0.4ms avg with 16 servers):

# What to choose?



Response Distribution M/G Process (10000 QPS @ 0.4ms avg with 16 servers):

Information Free "Static"

# What to choose?



Response Distribution M/G Process (10000 QPS @ 0.4ms avg with 16 servers):

"Dynamic" based on local knowledge

# What to choose?



Response Distribution M/G Process (10000 QPS @ 0.4ms avg with 16 servers):

"Dynamic" based on global knowledge

**What to choose?**

HAProxy recommends [least connections](#) as being strictly dominate to choice of 2 with an efficient impl

This matches the math and literature absent information.

Google allows servers to communicate back with clients to [adjust weights](#) in RR. Very clever.

# Stateful Load Balancing

State makes the
problem different

N

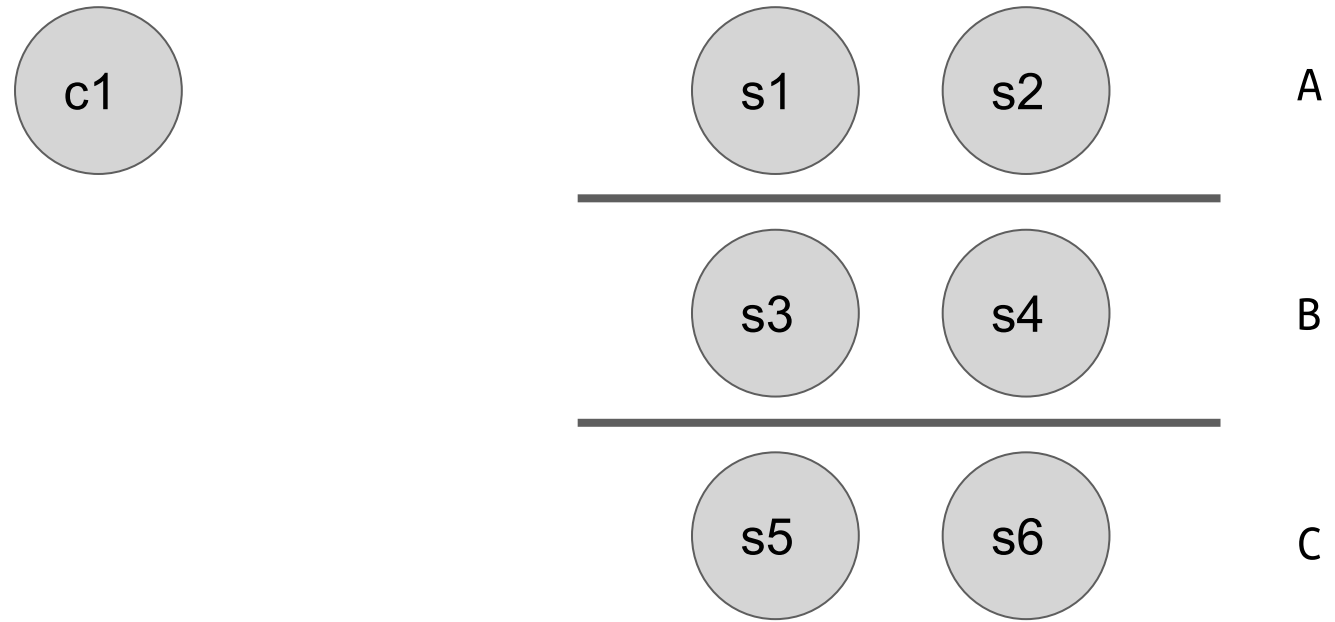**What makes datastores special?**

The node you hit matters!

- Postgres: master, replica

- ZooKeeper: leader, followers

- CockroachDB: lease holder

N

**What makes Cassandra special?**

1. For any piece of data we typically have one **replica** per availability zone

2. Depending on the **consistency** we may need to hop to more hosts

3. Datastores have hiccups frequently (drives mostly)

4. Our network latency is **asymmetric**

N

# Stateful load balancing with real networks



|     | A          | B          | C          |
|-----|------------|------------|------------|
| A   | **150us**  | 800us      | 250us      |
| B   | 800us      | **220us**  | 850us      |
| C   | 380us      | 700us      | **160us**  |

# DataStax Java Driver for Apache Cassandra®

**DataStax Java Driver 3.x for Apache Cassandra®**

**No Token?** Round Robin

**Token Aware?** Hash key, shuffle replicas*, return first. (**random subsetting**)

**DataStax Java Driver 3.x for Apache Cassandra®**

Slow to react to slow coordinators, erroring coordinators, paused coordinators, etc …

Traffic often goes **cross-zone**
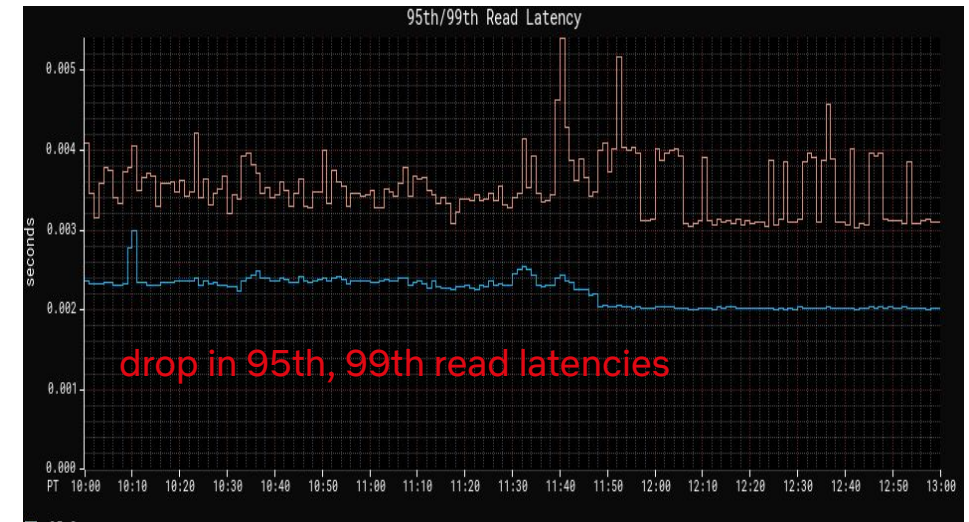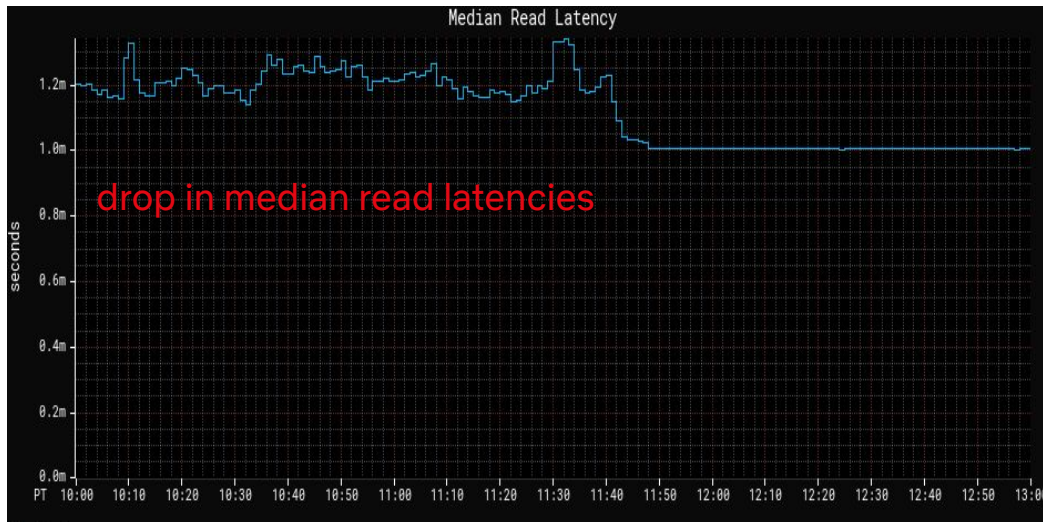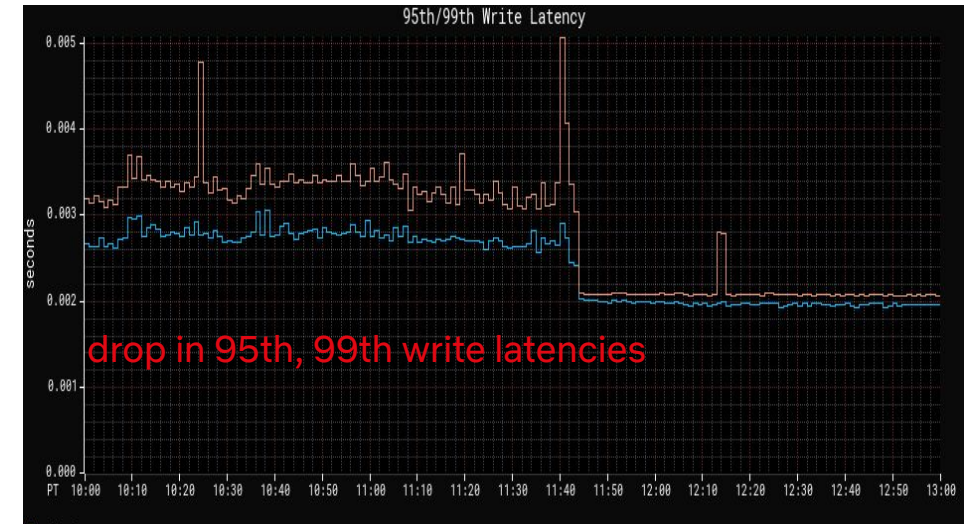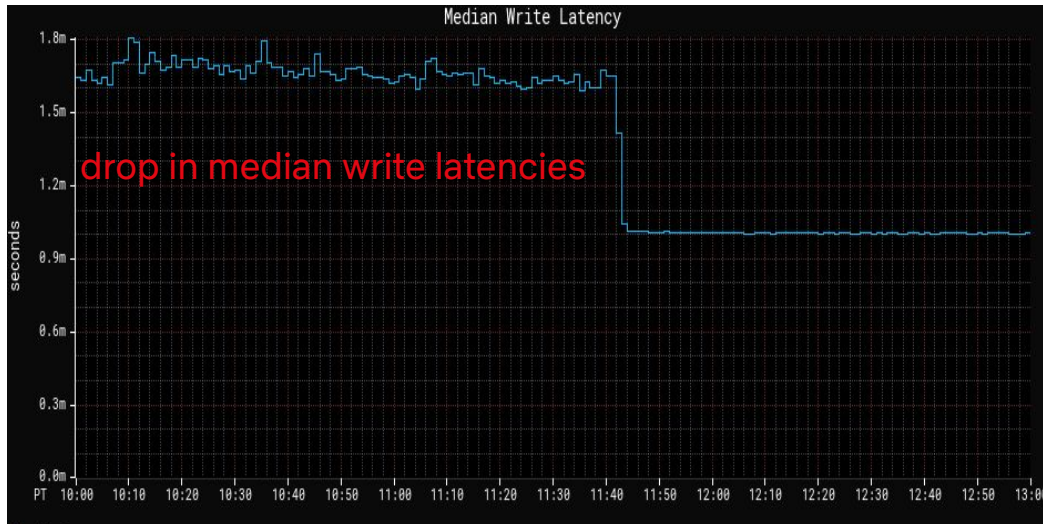
## No Token?
Round Robin

## Token Aware?
Hash key, shuffle replicas, return least loaded between first and second.

Avoids very slow replicas!

Basically choice of 2 over random subsets! Nice!

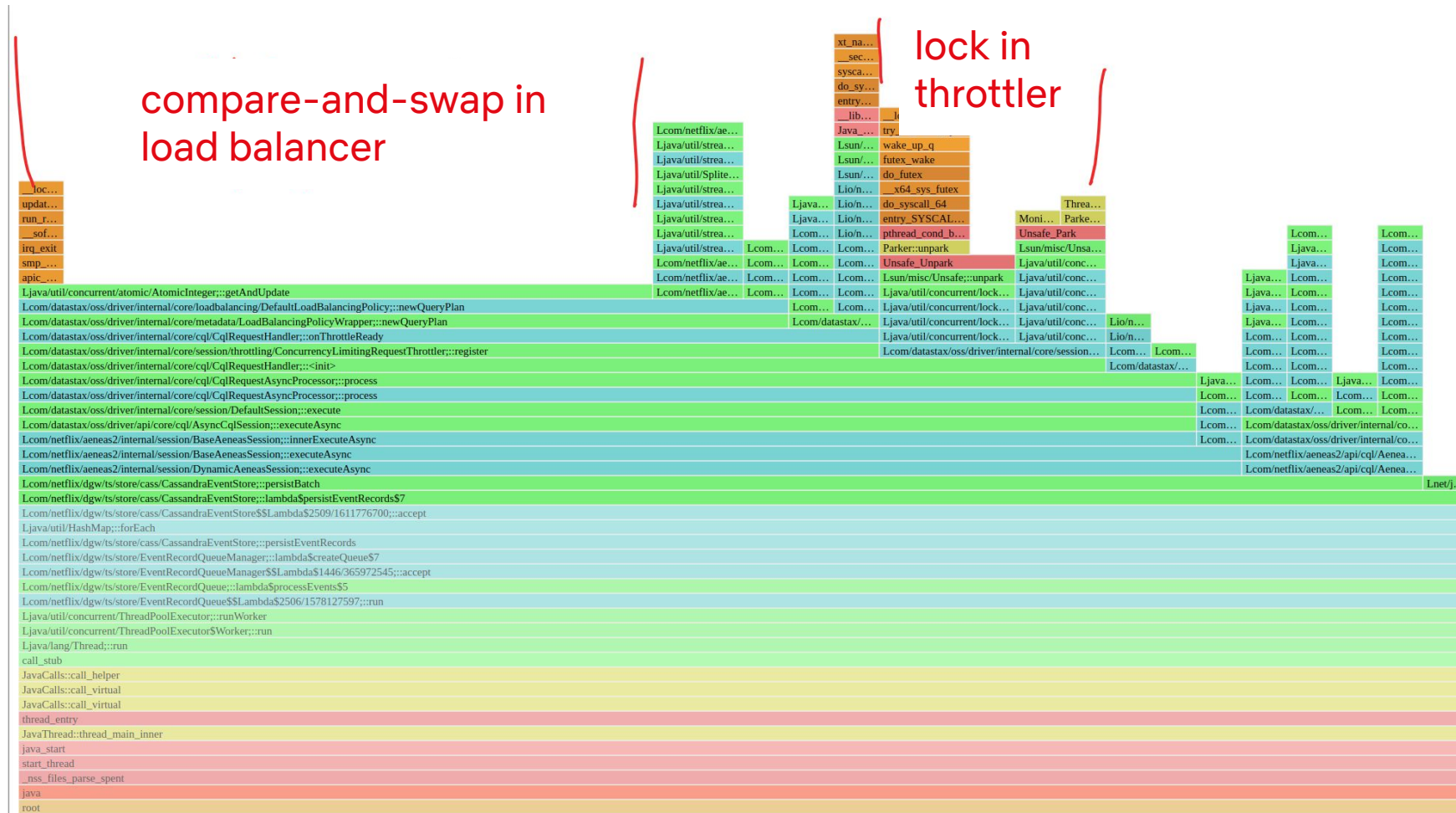# DataStax Java Driver 4.x for Apache Cassandra®

**DataStax Java Driver 4.x for Apache Cassandra®**

Perf regression with **high-throughput** cases

We needed to do 20k QPS per client to Cassandra and Datastax 4.x could barely do 8k.

# DataStax Java Driver 4.x for Apache Cassandra®

# Pays expensive **compare and update** and a **lock acquire-release**



compare-and-swap in load balancer

lock in throttler

**DataStax Java Driver 4.x for Apache Cassandra®**

# Pays expensive **compare and update** and a **lock acquire-release**

```java
LOG.trace("[{}] Prioritizing {} local replicas", logPrefix, replicaCount);

// Round-robin the remaining nodes
ArrayUtils.rotate(
    currentNodes,                                    DefaultLoadBalancingPolicy#newQueryPlan
    replicaCount,
    length: currentNodes.length - replicaCount,
    roundRobinAmount.getAndUpdate(INCREMENT));

QueryPlan plan = currentNodes.length == 0 ? QueryPlan.EMPTY : new SimpleQueryPlan(currentNodes);
return maybeAddDcFailover(request, plan);
}
```

# Weighted Least Loaded

Started with fixing compare-and-swap, ended up rewriting the algorithm

**No Token?**

Chose 8 random nodes

**Token Aware?**

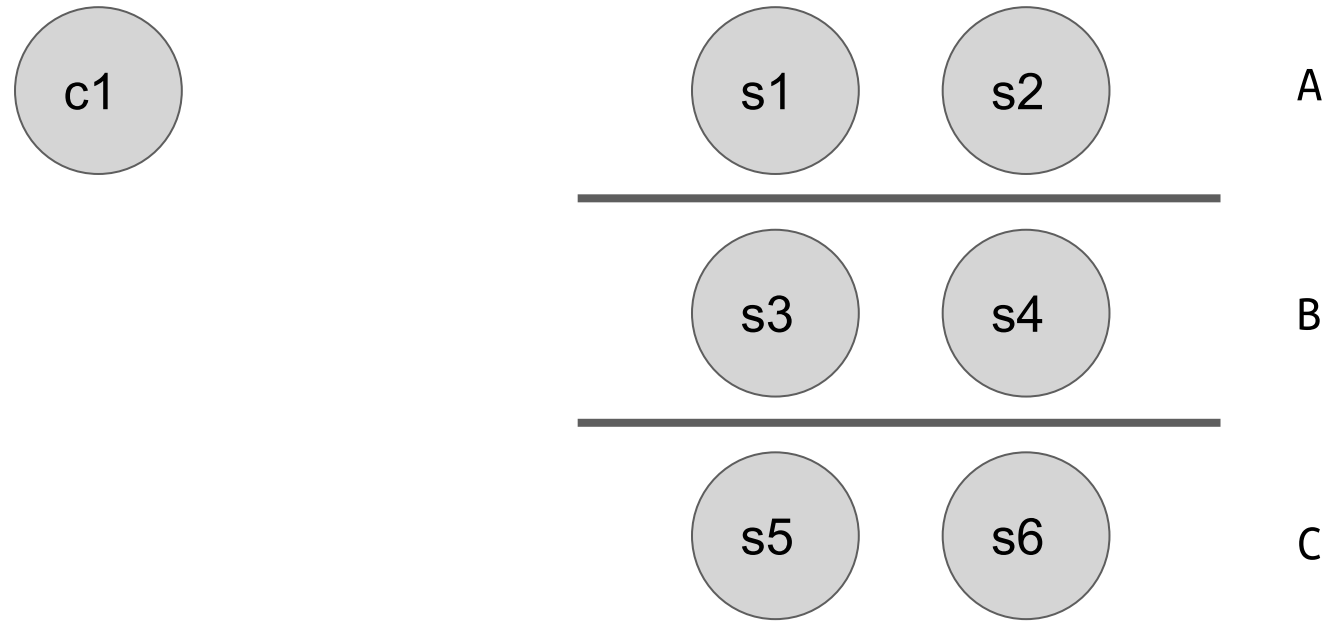Choose all RF replicas and 8-RF random

**Weight concurrency by:**
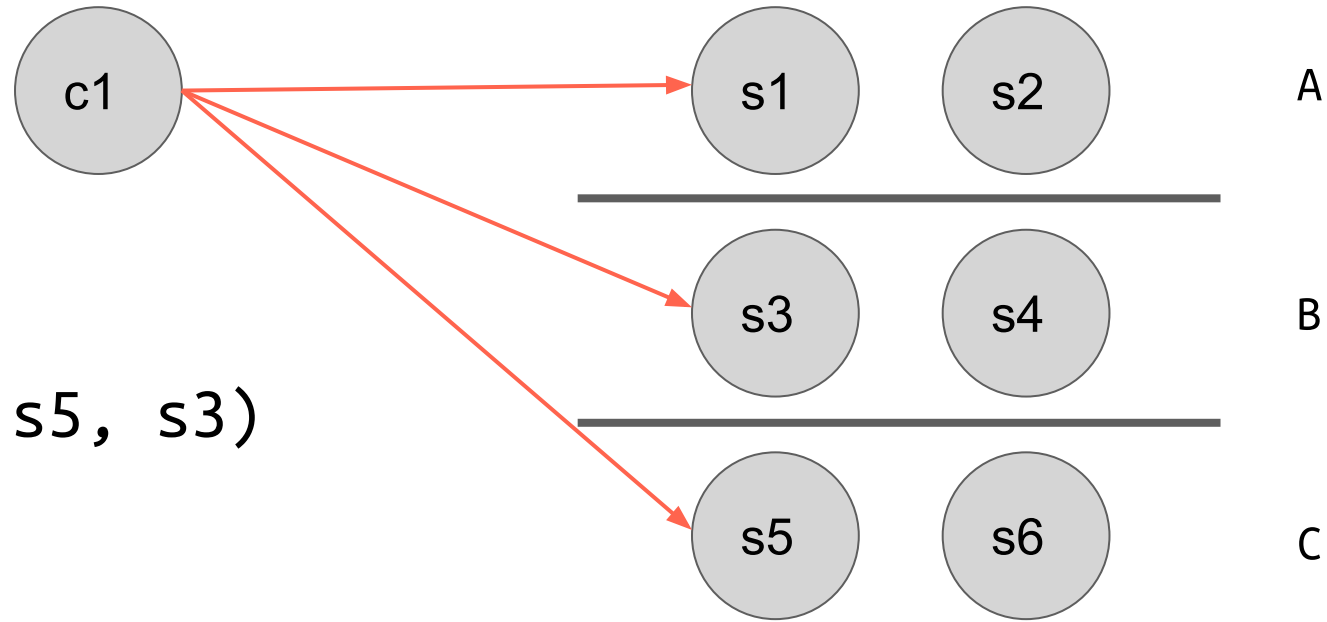
!Rack = 4

!Replica = 12

Unhealthy = 64

Sort the sublist. Done!

# Stateful load balancing with real networks

c1

s1    s2    A

s3    s4    B

s5    s6    C

|   | A | B | C |
|---|---|---|---|
| A | **150us** | 800us | 250us |
| B | 800us | **220us** | 850us |
| C | 380us | 700us | **160us** |

**LOCAL_ONE
(Control)**

c1

s1  s2    A

s3  s4    B

s5  s6    C

set(x=0)
replicas(x) = (s1, s5, s3)

End to End Latency = Latency (L) + Processing (R)

E_LO = ⅓ (L(A, A) + R) + ⅓ (L(A, B) + R) + ⅓ (L(A, C) + R)
Let R = 100us
E_LO = ⅓ (150 + 100) + ⅓ (800 + 100) + ⅓ (250 + 100) = **500us**

**LOCAL_ONE
(WLLLB)**



set(x=0)
replicas(x) = (s1, s5, s3)

End to End Latency = Latency (L) + Processing (R)

E_LO = L(A, A) + R
Let R = 100us
E_LO = 150 + 100 = **250us** (50% reduction)

**LOCAL_QUORUM
(Control)**



set(x=0)
replicas(x) = (s1, s5, s3)

E_LQ = ⅓ (**L(A, A)** + min(R, L(A, C) + R))
⅓ (**L(A, B)** + min(R, L(B, A) + R))
⅓ (**L(A, C)** + min(R, L(C, A) + R))

Let R = 100us
E_LQ = ⅓ (150 + 350) + ⅓ (800 + 900) + ⅓ (250 + 480) = **980us**

**LOCAL_QUORUM
(Control)**



set(x=0)
replicas(x) = (s1, s5, s3)

$$E\_LQ = \tfrac{1}{3} \; (L(A, A) + \mathbf{min}(R, L(A, C) + R))$$
$$\tfrac{1}{3} \; (L(A, B) + \mathbf{min}(R, L(B, A) + R))$$
$$\tfrac{1}{3} \; (L(A, C) + \mathbf{min}(R, L(C, A) + R))$$

Let R = 100us
$$E\_LQ = \tfrac{1}{3} \; (150 + 350) + \tfrac{1}{3} \; (800 + 900) + \tfrac{1}{3} \; (250 + 480) = \mathbf{980us}$$

**LOCAL_QUORUM (WLLLB)**

```
set(x=0)
replicas(x) = (s1, s5, s3)
```

$$E\_LQ = \mathbf{L(A, A)} + \min(R, L(A, C) + R)$$

```
Let R = 100us
E_LQ = 150 + 100 + 250 = 500us (50% reduction)
```

**LOCAL_QUORUM (WLLLB)**

set(x=0)
replicas(x) = (s1, s5, s3)

E_LQ = L(A, A) + **min(R, L(A, C) + R)**

Let R = 100us
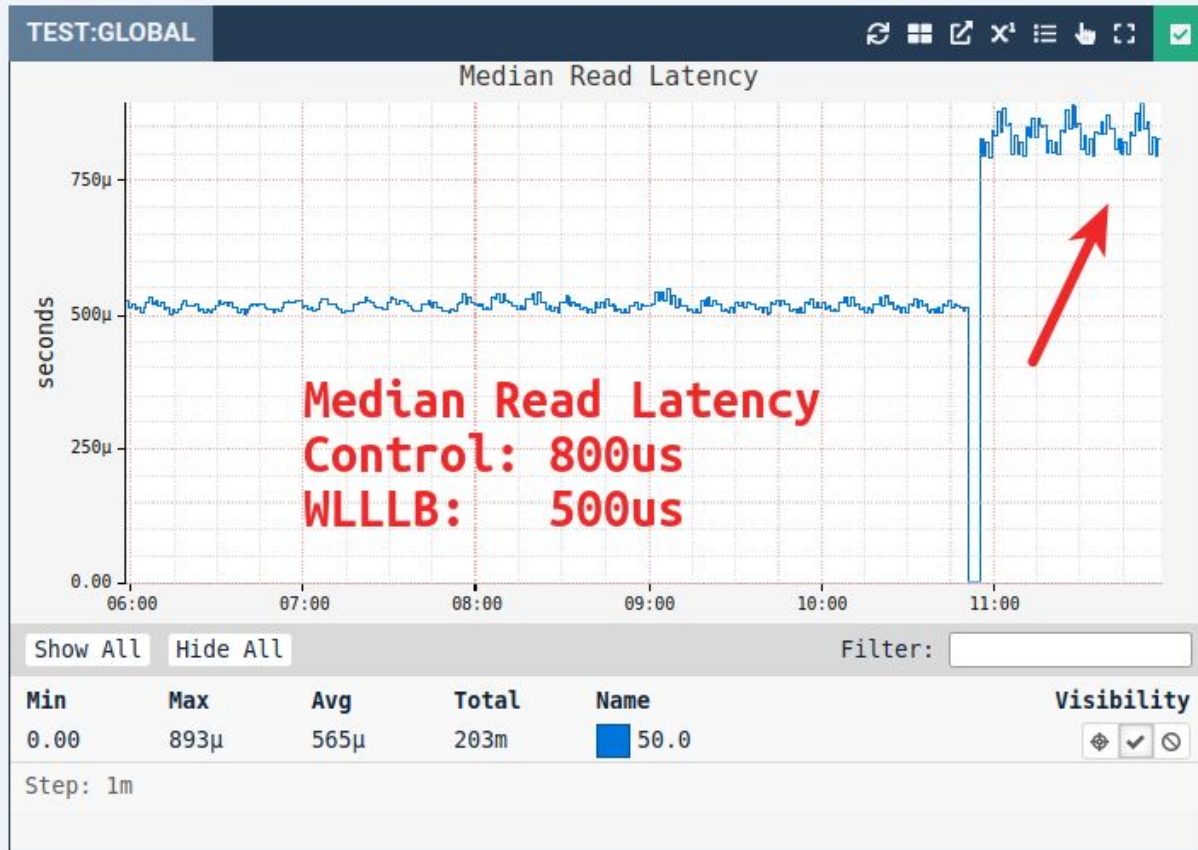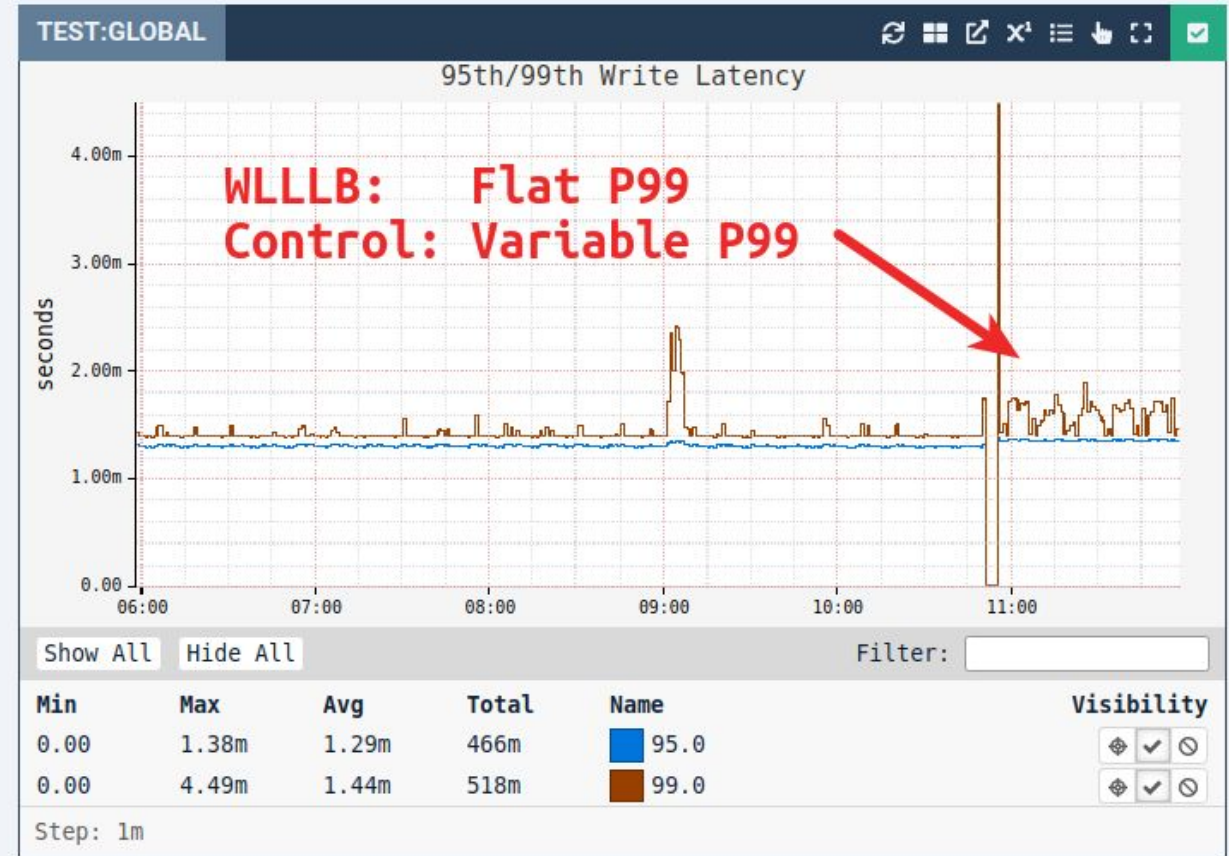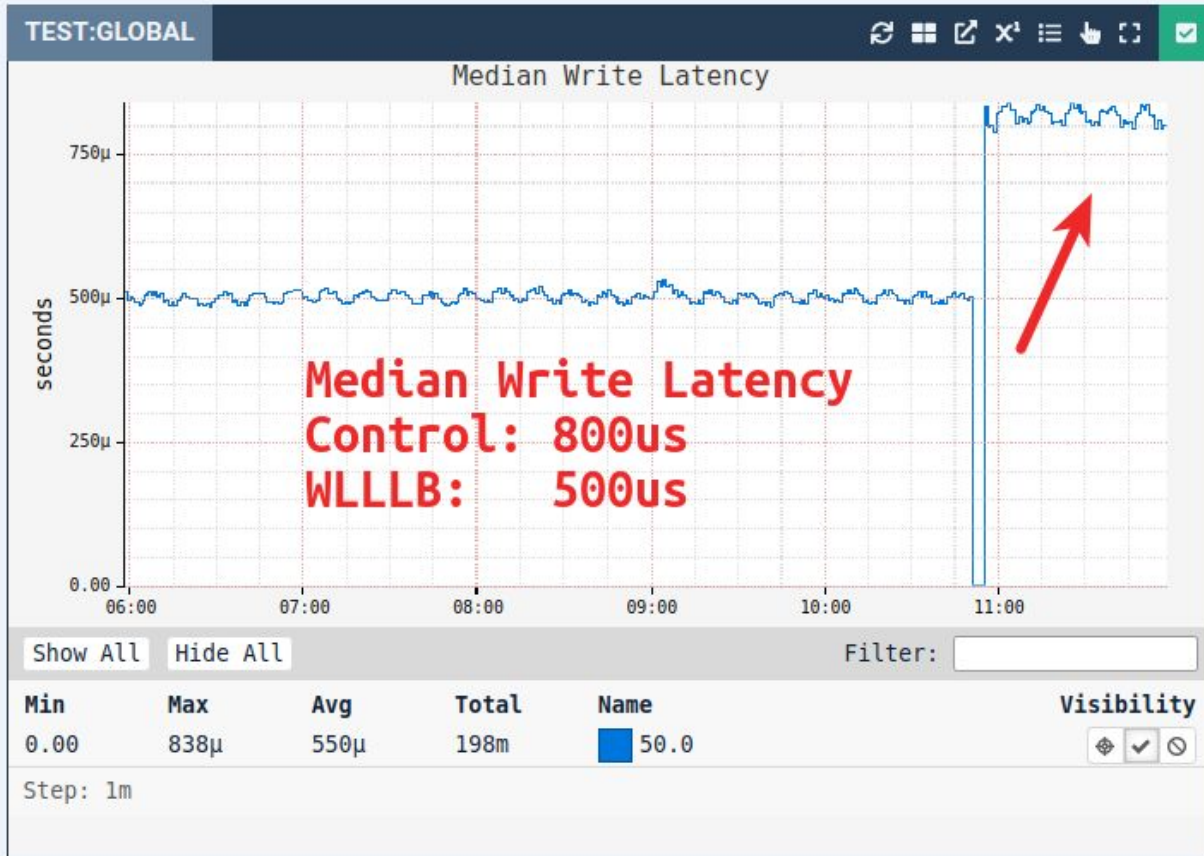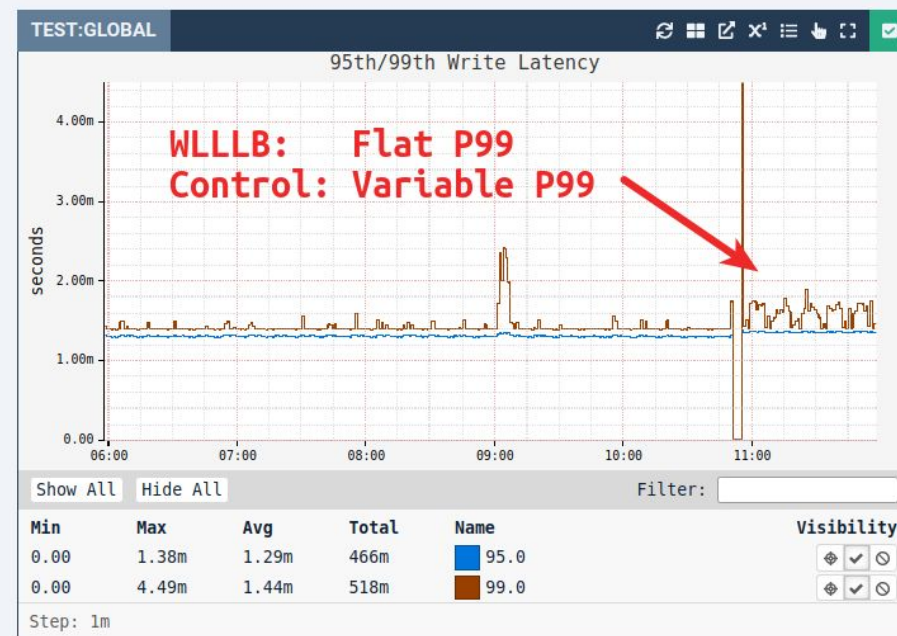E_LQ = 150 + 100 + 250 = 500us **(50% reduction)**

# Experiments

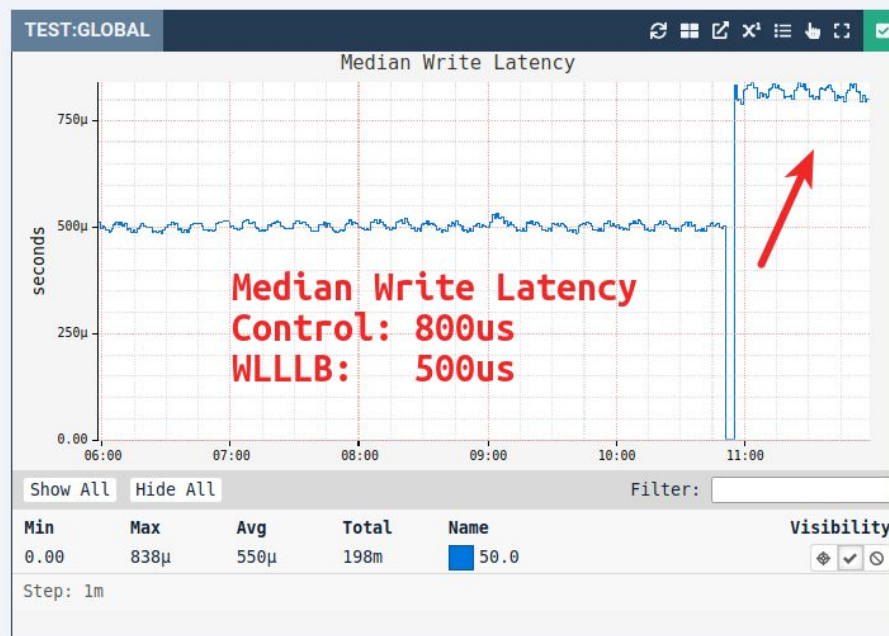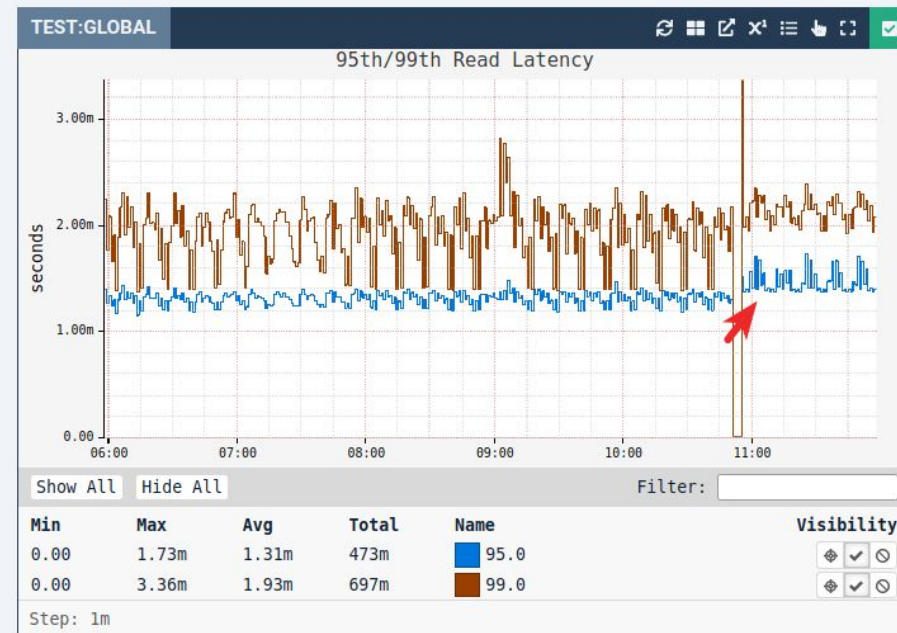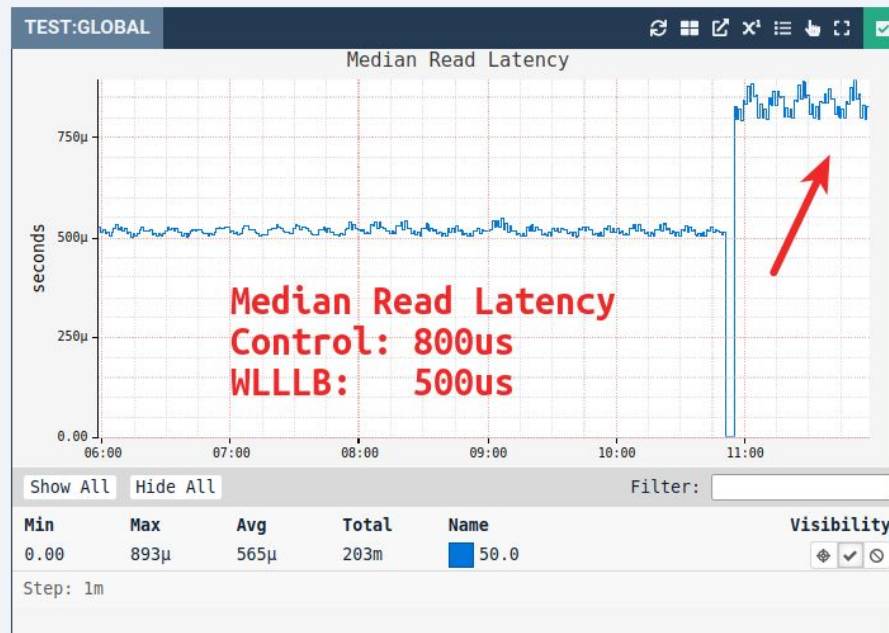# Synthetic Traffic

Apply Load
Measure Results

# Latency results
# LOCAL_ONE

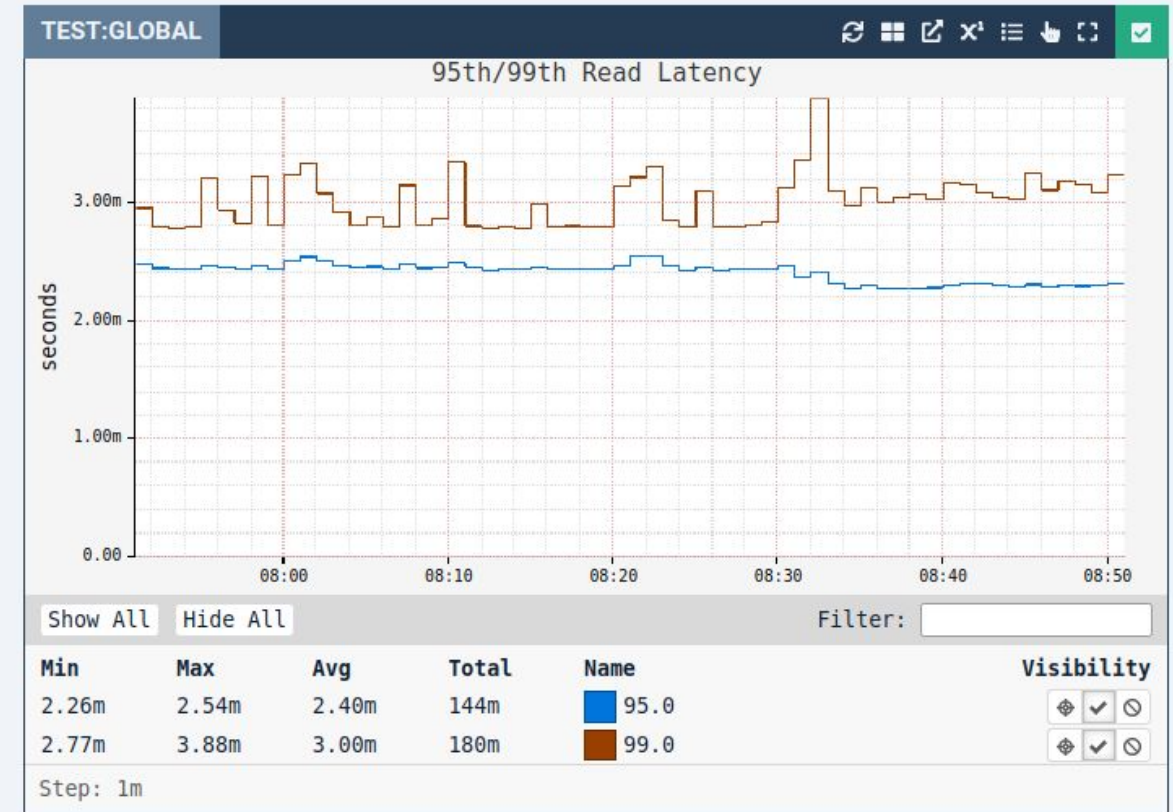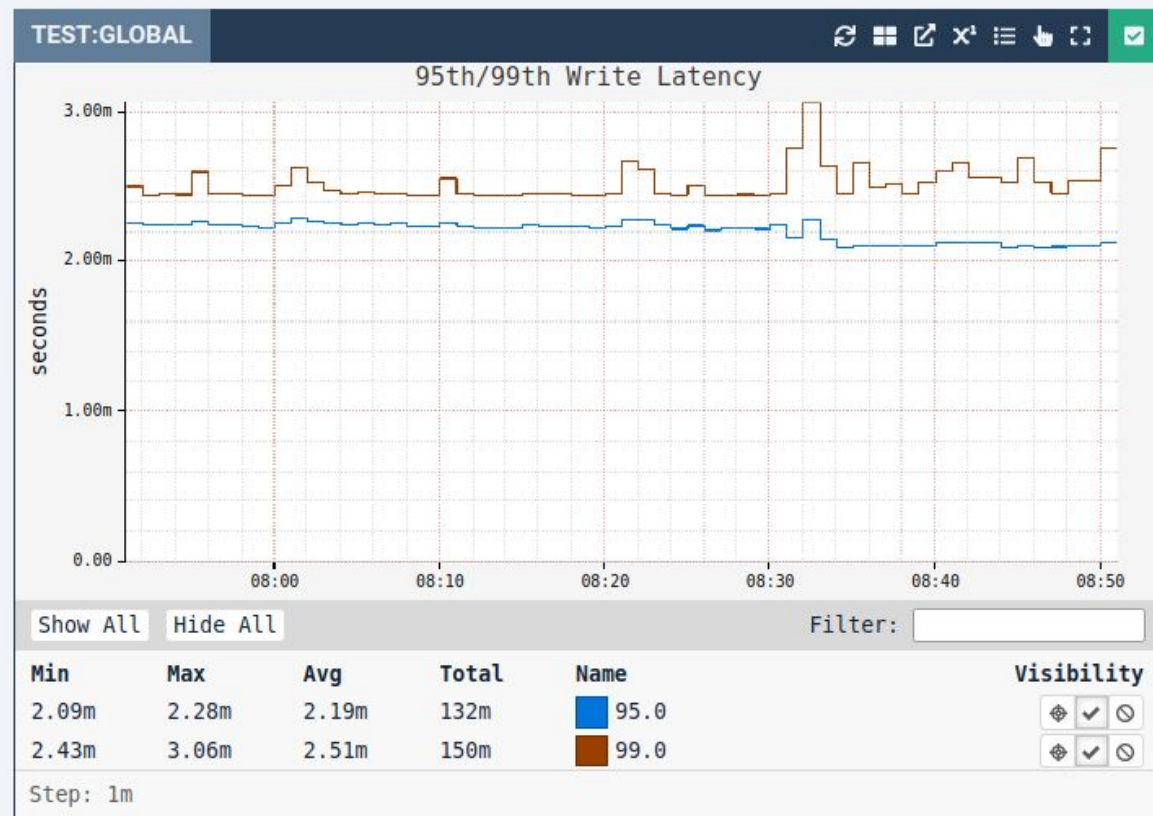# Latency results
# LOCAL_ONE

## Latency results LOCAL_ONE

## About a 40% improvement

# Latency results
# LOCAL_QUORUM



Read Latency

**TEST:GLOBAL**

Median Read Latency

WLLLB    Control

| Min | Max | Avg | Total | Name | | Visibility |
|-----|-----|-----|-------|------|--|------------|
| 1.32m | 1.55m | 1.38m | 83.0m | 🟦 50.0 | | |

Step: 1m

**TEST:GLOBAL**

95th/99th Read Latency

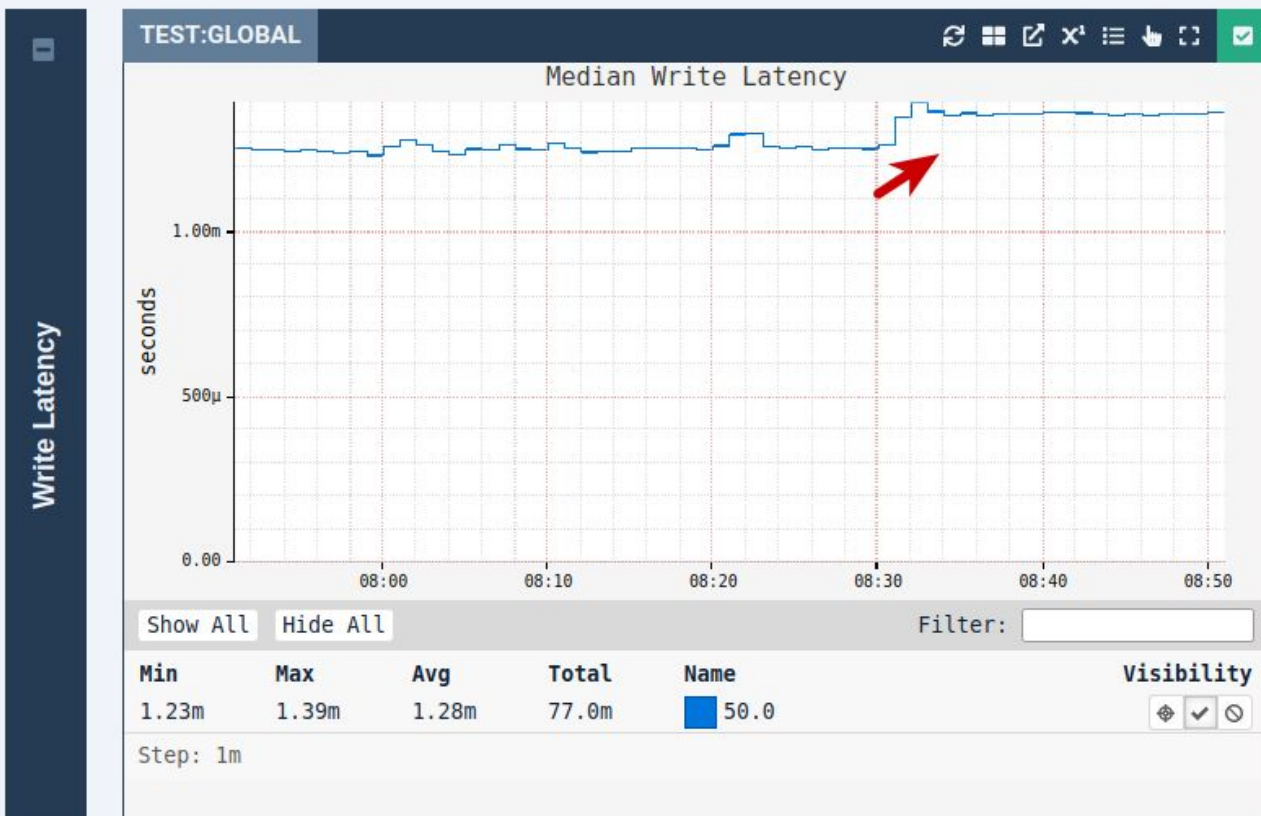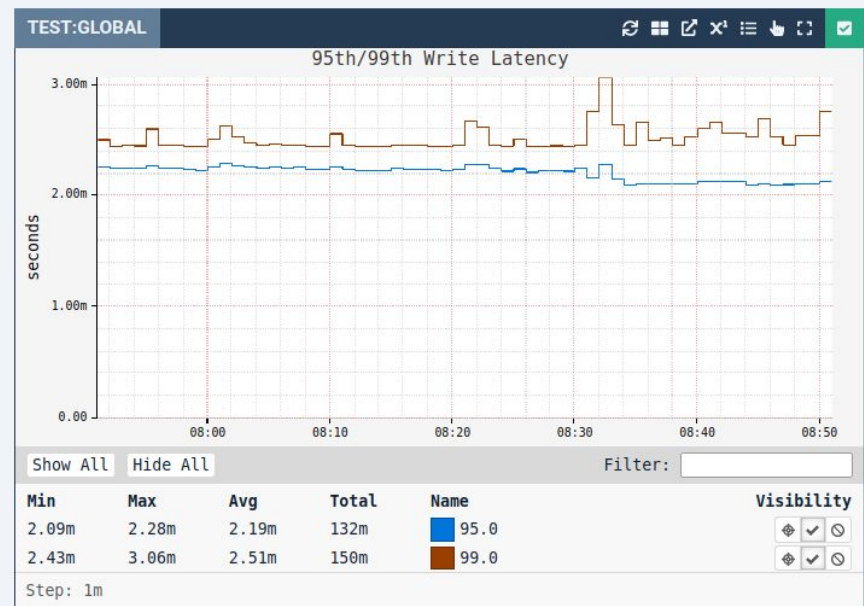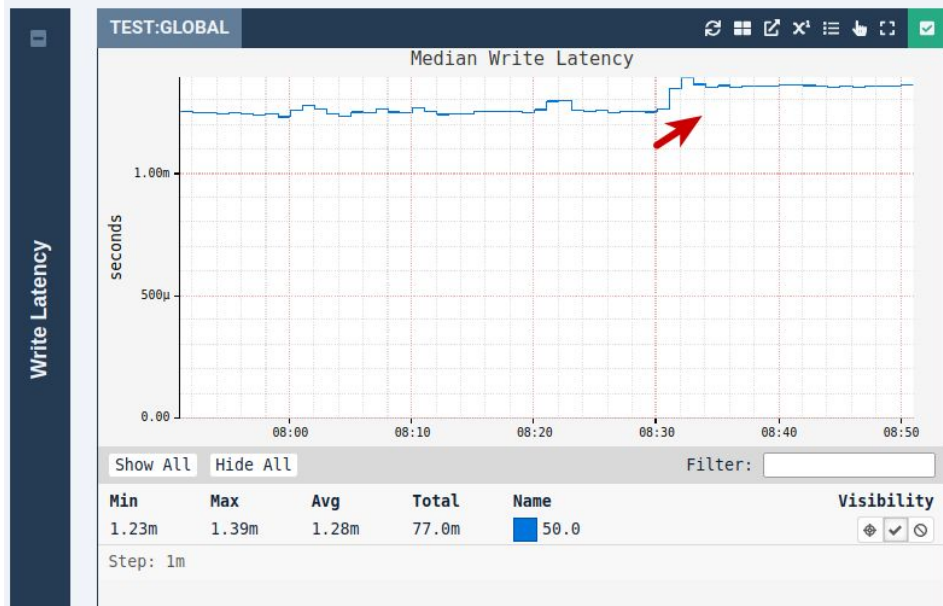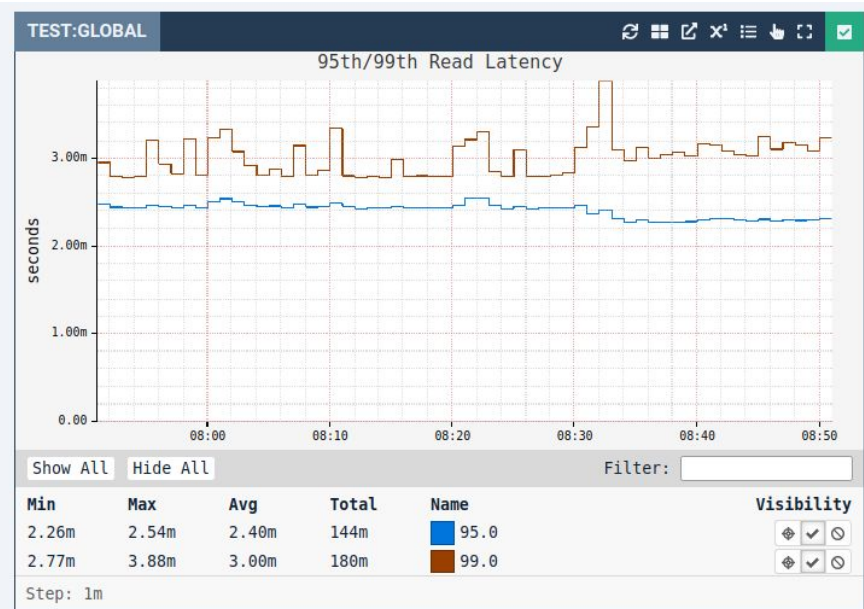| Min | Max | Avg | Total | Name | | Visibility |
|-----|-----|-----|-------|------|--|------------|
| 2.26m | 2.54m | 2.40m | 144m | 🟦 95.0 | | |
| 2.77m | 3.88m | 3.00m | 180m | 🟫 99.0 | | |

Step: 1m

# Latency results
# LOCAL_QUORUM

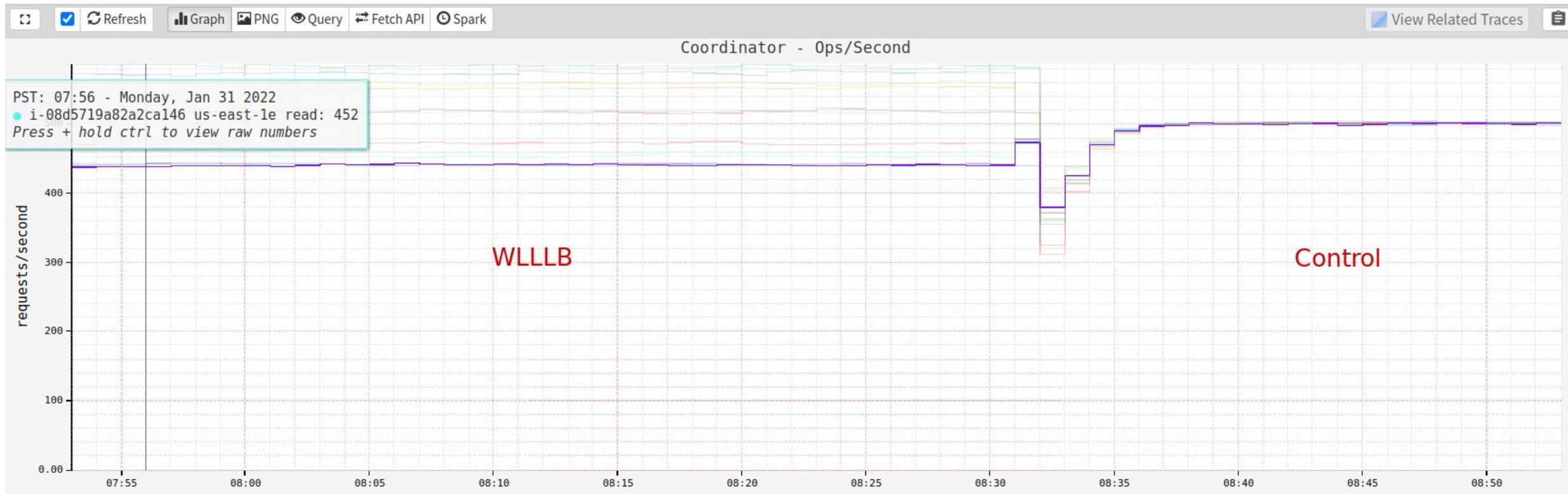# Latency results LOCAL_QUORUM

# About a 10% improvement

# Latency results

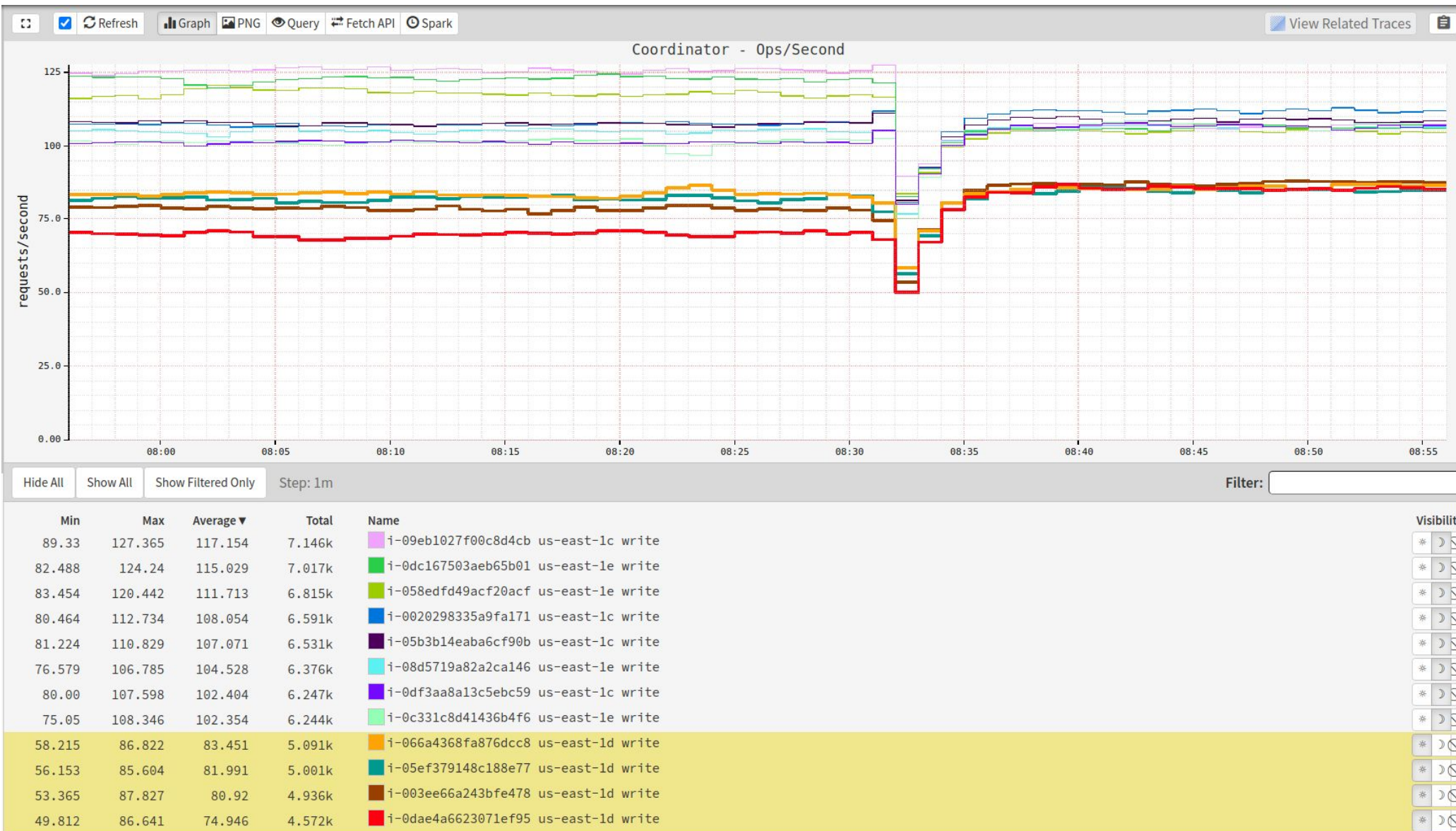| | WLLLB P50/P95/P99 Read (ms) | WLLLB P50/P95/P99 Write (ms) | Control P50/P95/P99 Read (ms) | Control P50/P95/P99 Write (ms) | Read Latency Difference | Write Latency Difference |
|---|---|---|---|---|---|---|
| LO-1 | 0.52/1.30/1.92 | 0.50/1.30/1.41 | 0.84/1.45/2.14 | 0.82/1.35/1.59 | **38%**/10%/10% | **39%**/4%/11% |
| LQ-1 | 1.33/2.42/2.90 | 1.21/2.15/2.45 | 1.52/2.25/3.07 | 1.36/2.06/2.48 | **12.5/**-7.5%/5.6% | **11%/**-4.3%/1.2% |
| LQ-2 | 1.40/2.56/4.45 | 1.27/2.08/2.46 | 1.55/2.32/3.93 | 1.32/2.03/2.47 | **10%**/-10%/-13% | **4%**/-5%/-1% |

Why the slight P95 regression in LQ? Theories:

1. Load Imbalance due to asymmetric latency
2. Dynamic Endpoint Snitch

# Load imbalance
# Reads

# Load imbalance
# Writes

# Network Delay

Force packet delay

Measure results

# Linux Traffic Control (tc)!

```
$ sudo tc qdisc show dev eth0
qdisc mq 8005: root
qdisc fq 0: parent 8005:4 limit 10000p flow_limit 100p buckets 1024 orphan_mask 1023 quantum 18030
initial_quantum 90150 low_rate_threshold 550Kbit refill_delay 40.0ms
qdisc fq 0: parent 8005:3 limit 10000p flow_limit 100p buckets 1024 orphan_mask 1023 quantum 18030
initial_quantum 90150 low_rate_threshold 550Kbit refill_delay 40.0ms
qdisc fq 0: parent 8005:2 limit 10000p flow_limit 100p buckets 1024 orphan_mask 1023 quantum 18030
initial_quantum 90150 low_rate_threshold 550Kbit refill_delay 40.0ms
qdisc fq 0: parent 8005:1 limit 10000p flow_limit 100p buckets 1024 orphan_mask 1023 quantum 18030
initial_quantum 90150 low_rate_threshold 550Kbit refill_delay 40.0ms
```

# Netem to the rescue
## ([tc-netem](tc-netem))

```
# Server adds 10ms delay
server$ sudo tc qdisc replace dev eth0 root netem delay 10ms

# Client now observes 10ms additional latency on all requests
client$ ping 100...
…
64 bytes from 100...: icmp_seq=525 ttl=64 time=0.215 ms
64 bytes from 100...: icmp_seq=526 ttl=64 time=0.212 ms
# When netem was enabled
64 bytes from 100...: icmp_seq=527 ttl=64 time=10.2 ms
64 bytes from 100...: icmp_seq=528 ttl=64 time=10.2 ms
64 bytes from 100...: icmp_seq=529 ttl=64 time=10.2 ms
64 bytes from 100...: icmp_seq=530 ttl=64 time=10.2 ms


# Now Revert on server
server$ sudo tc qdisc replace dev eth0 root mq
```

# Netem to the rescue
## ([tc-netem](tc-netem))

```
# You can also use netem to simulate packet loss, corruption,
duplication, reordering and other TCP issues.
# For example you could add a distribution of delay with


$ tc qdisc change dev eth0 root netem delay 10ms 4ms distribution normal
```
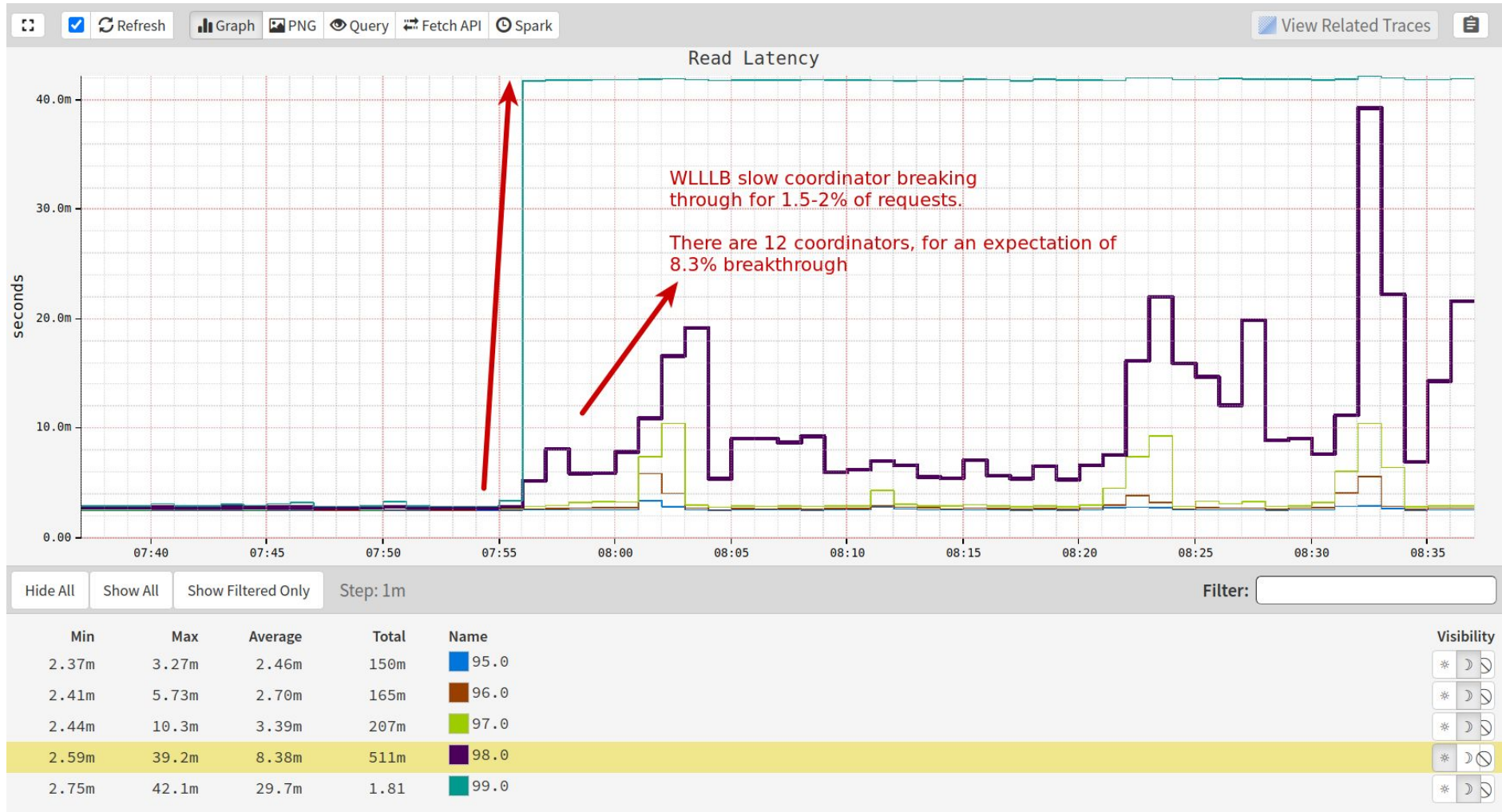
# Slow coordinators



Coordinator - Ops / Node

WLLLB shedding load off the affected coordinator

| Min | Max | Average | Total | Name | | Visibility |
|-----|-----|---------|-------|------|--|------------|
| 132 | 603 | 345 | 20.4k | ■ i-0da▒▒▒▒▒ us-east-1d | | ☀ ☽ ◖ |

# Slow coordinators

## Limited latency impact in ⅔ zones

# Slow coordinators

**1/12 = 8.3% should have been affected**

**But only 1.5% were**

# Garbage Collection

Simulate pauses

Measure results

# STOP + CONT

```
# pause.sh
while [ 1 ]
do
sudo -u www-data kill -STOP $(pgrep -f CassandraDaemon)
# Duration of pause
sleep 20
sudo -u www-data kill -CONT $(pgrep -f CassandraDaemon)
# Interval between pauses
sleep 30
done
```

# Slow coordinators

## Simulate "GC" pause via stopping the Java process.

# Slow coordinators

**1/12 = 8.3% should have been affected**

**But only .1% were**

# Real World Results

Apply Real Load
Measure Results

**Watch Graphs Drop**

# Service #1 - LOCAL_ONE



P50 1.1ms -> 0.7ms = **36%** improvement
P95 1.9ms -> 1.4ms = **26%** improvement
Local One workload

# Service #1 - LOCAL_ONE
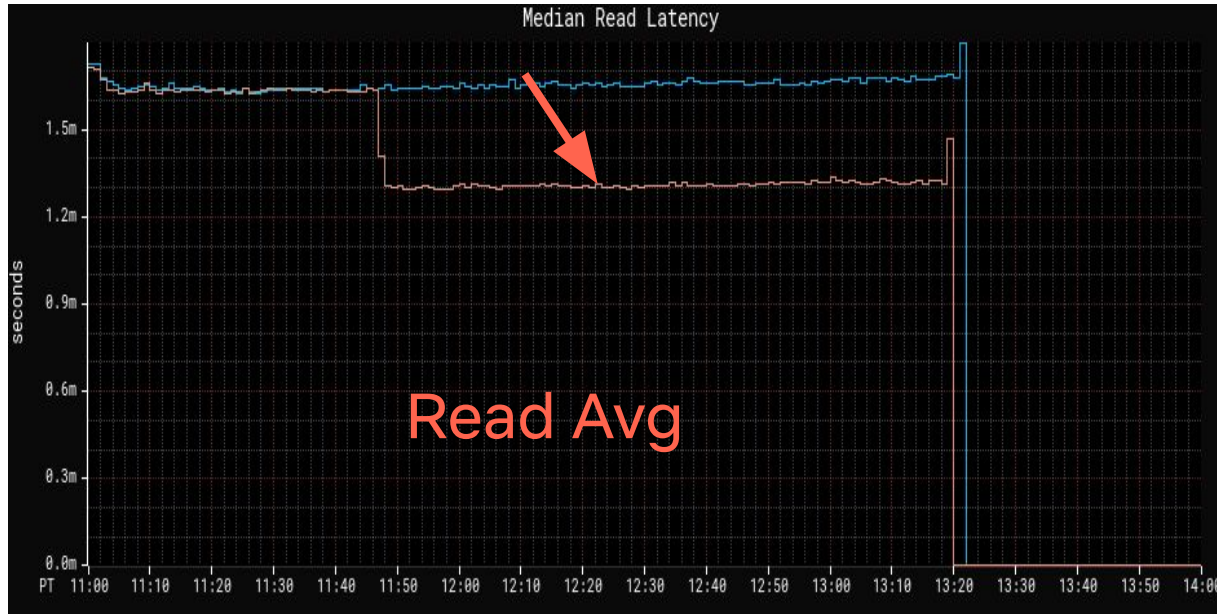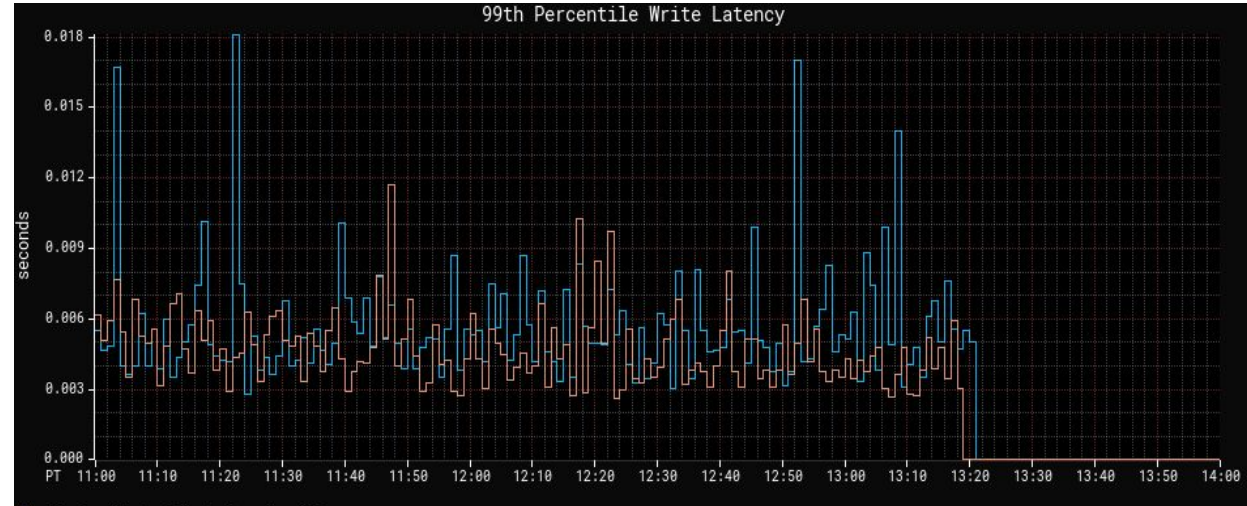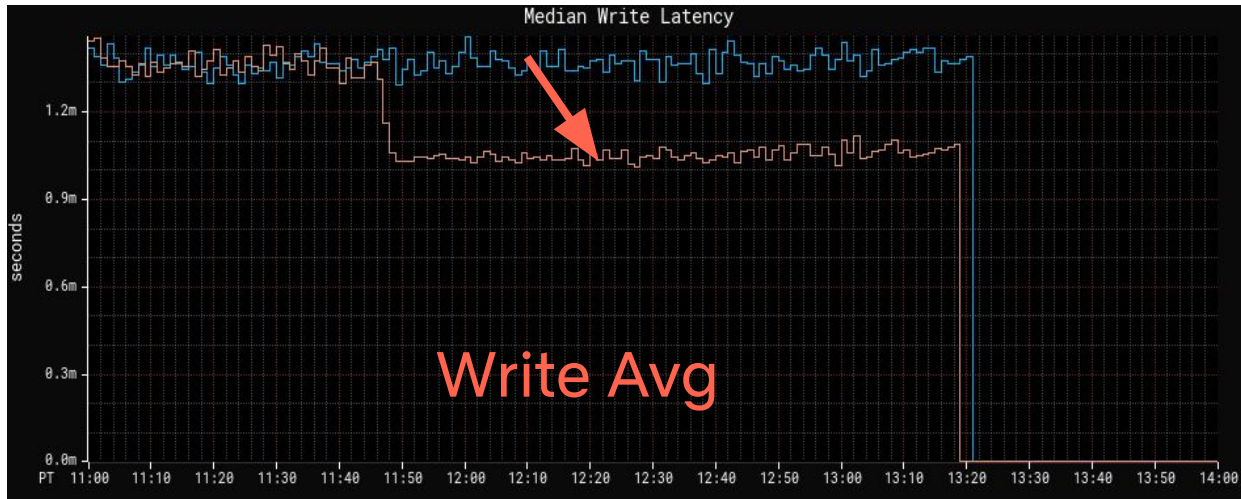


P50 1.2ms -> 0.7ms = **41%** improvement

P95 2.2ms -> 1.7ms = **22%** improvement

Local One workload

# Service #2 - LOCAL_QUORUM



P50 2.0ms -> 1.6ms = **20%** improvement

P95 2.8ms -> 2.2ms = **22%** improvement

LWT (Local Serial) workload

# Service #3 - LOCAL_ONE



P50 1.6ms -> 1.2ms = **25%** improvement

P99 5.0ms -> 4.2ms = **16%** improvement
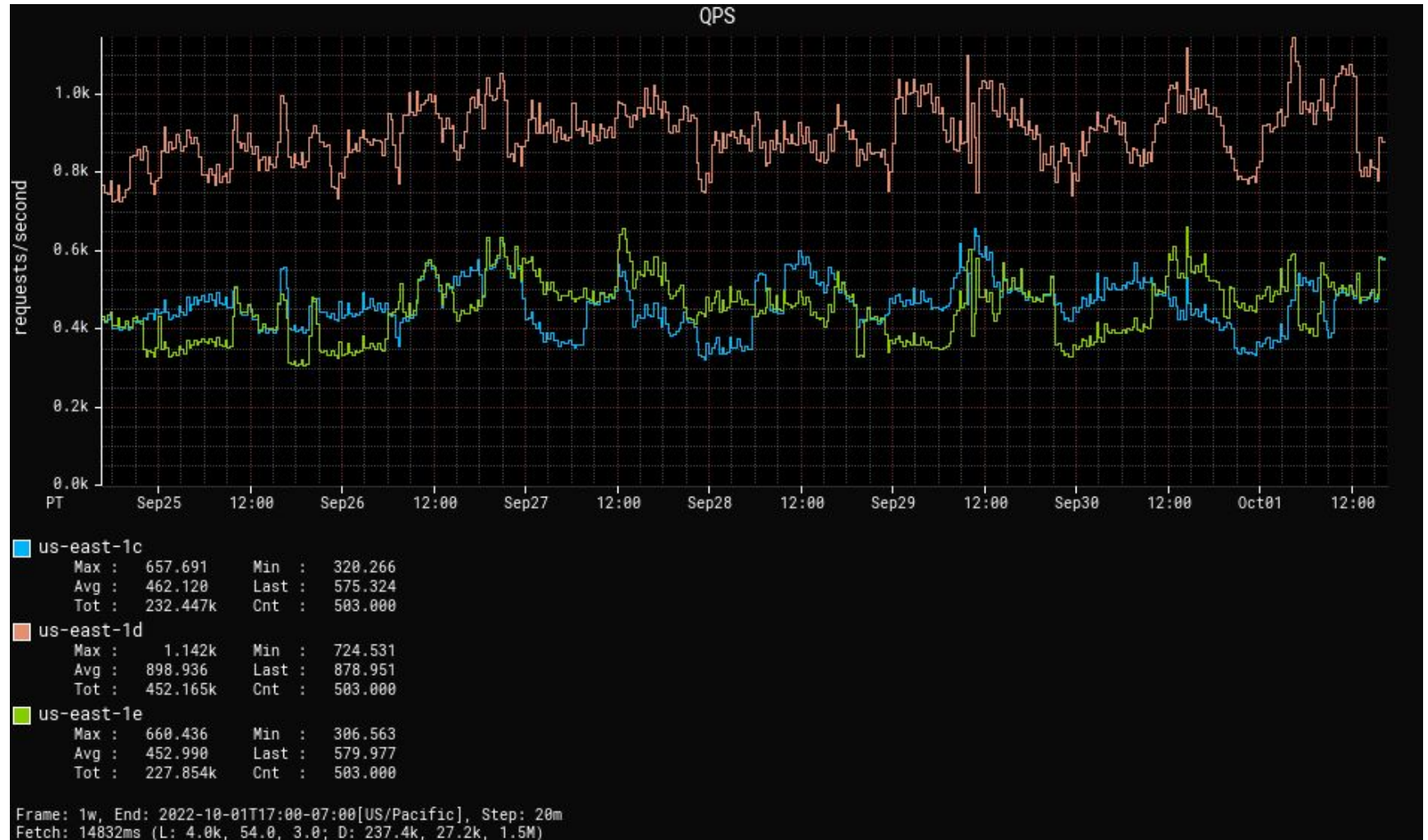
Local one workload

# Service #3 - LOCAL_ONE



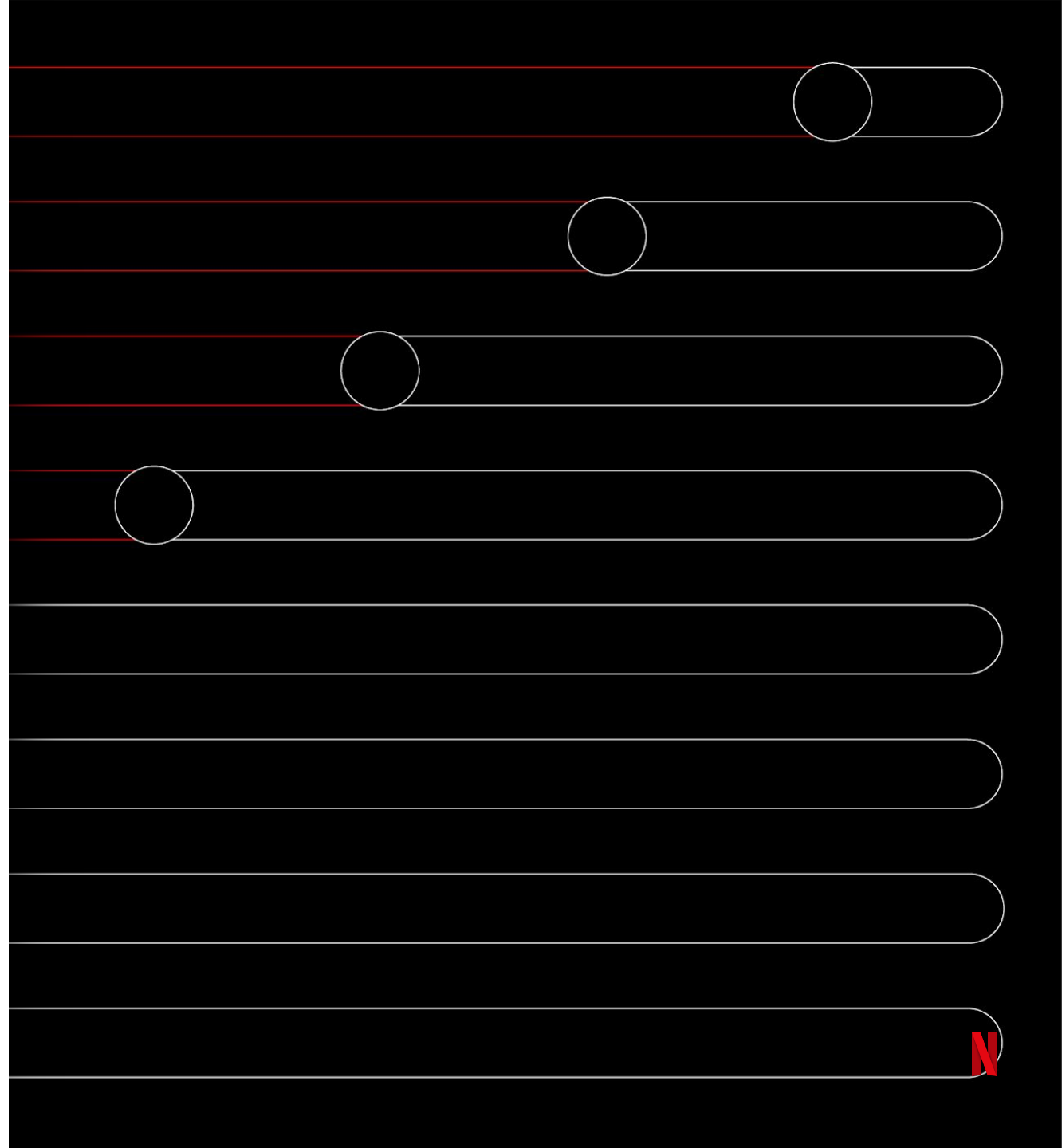P50 1.3ms -> 0.9ms = **31%** improvement

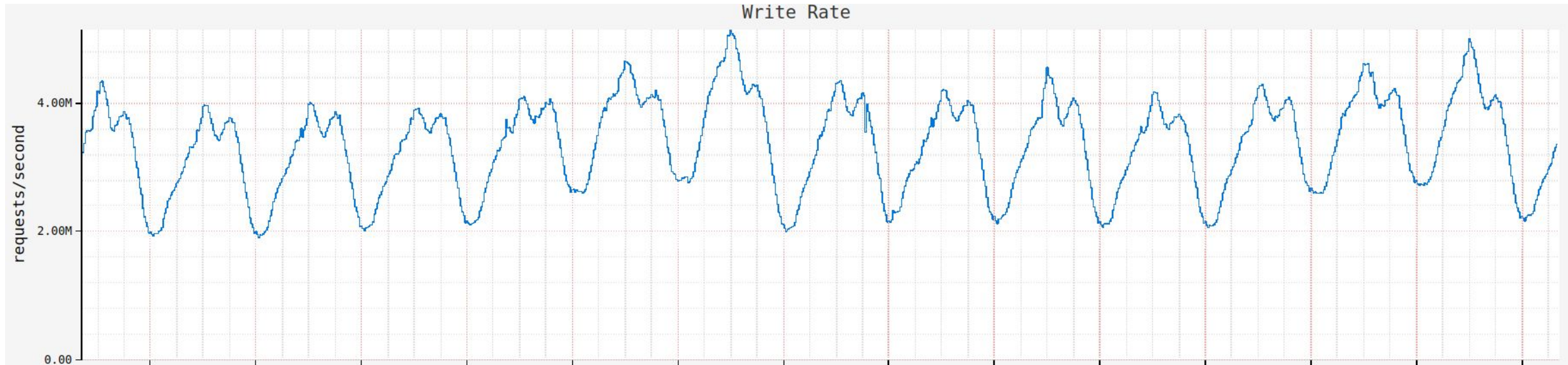P99 6.0ms -> 6.0ms = **~0%** improvement

Local one workload
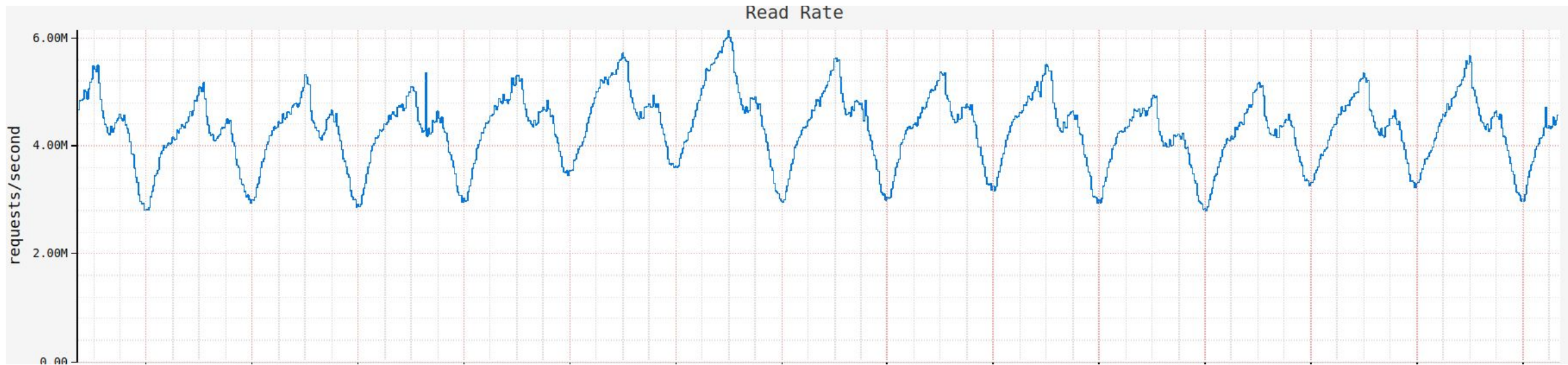
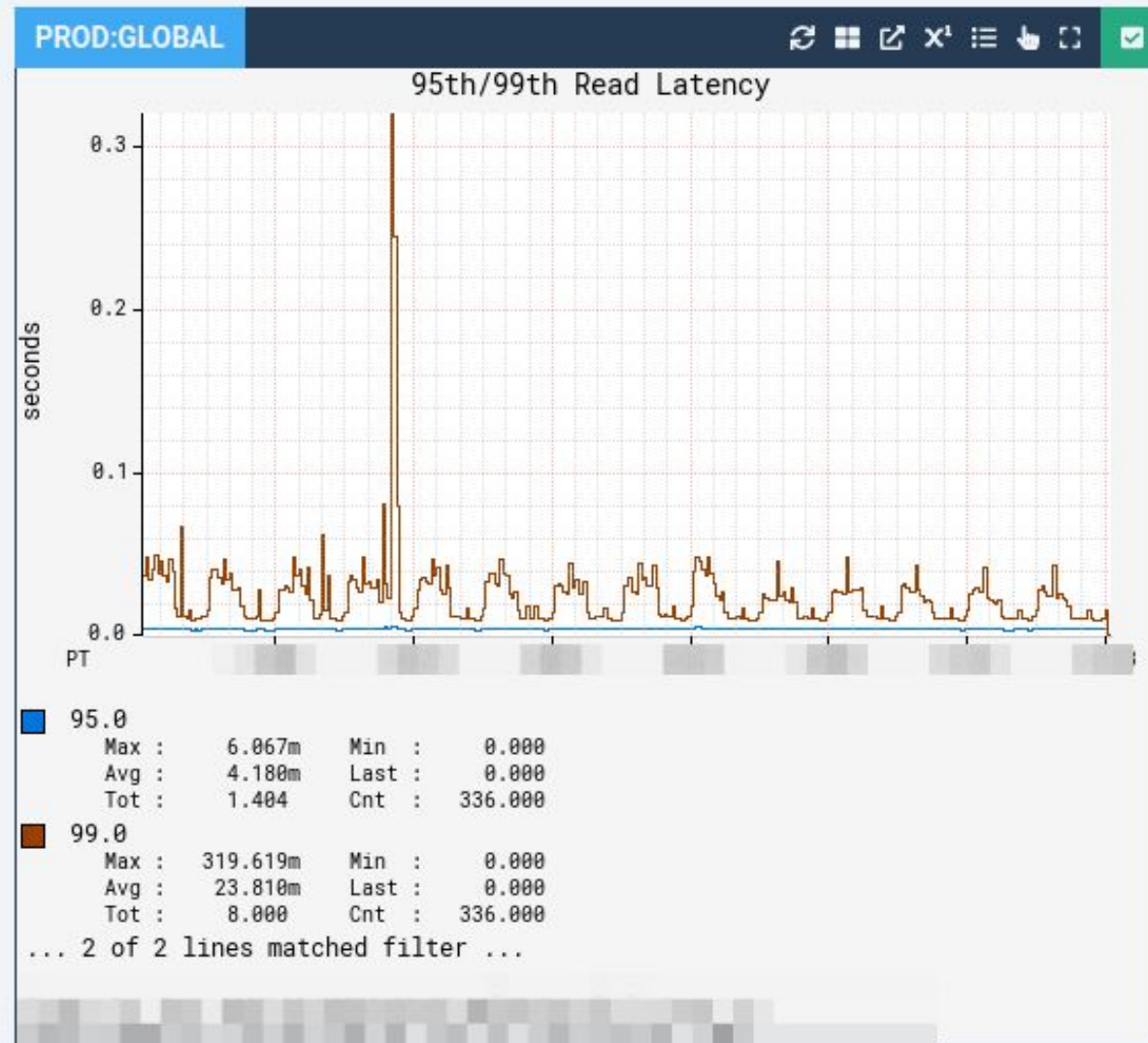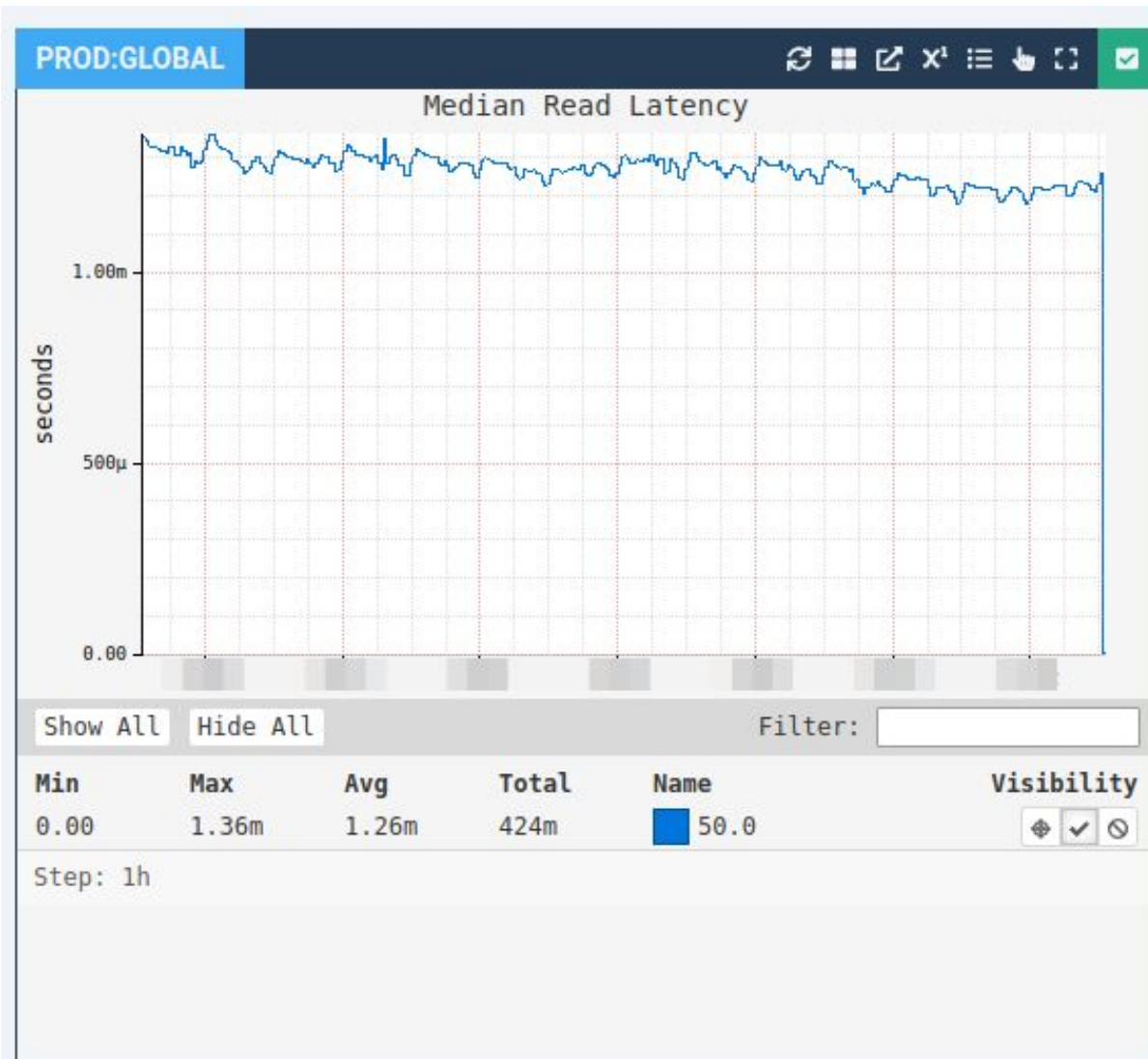# Uneven distribution of requests across zones

# At Scale?

# Scale?


Write Rate

Peak Traffic is 5 Million Writes per Second
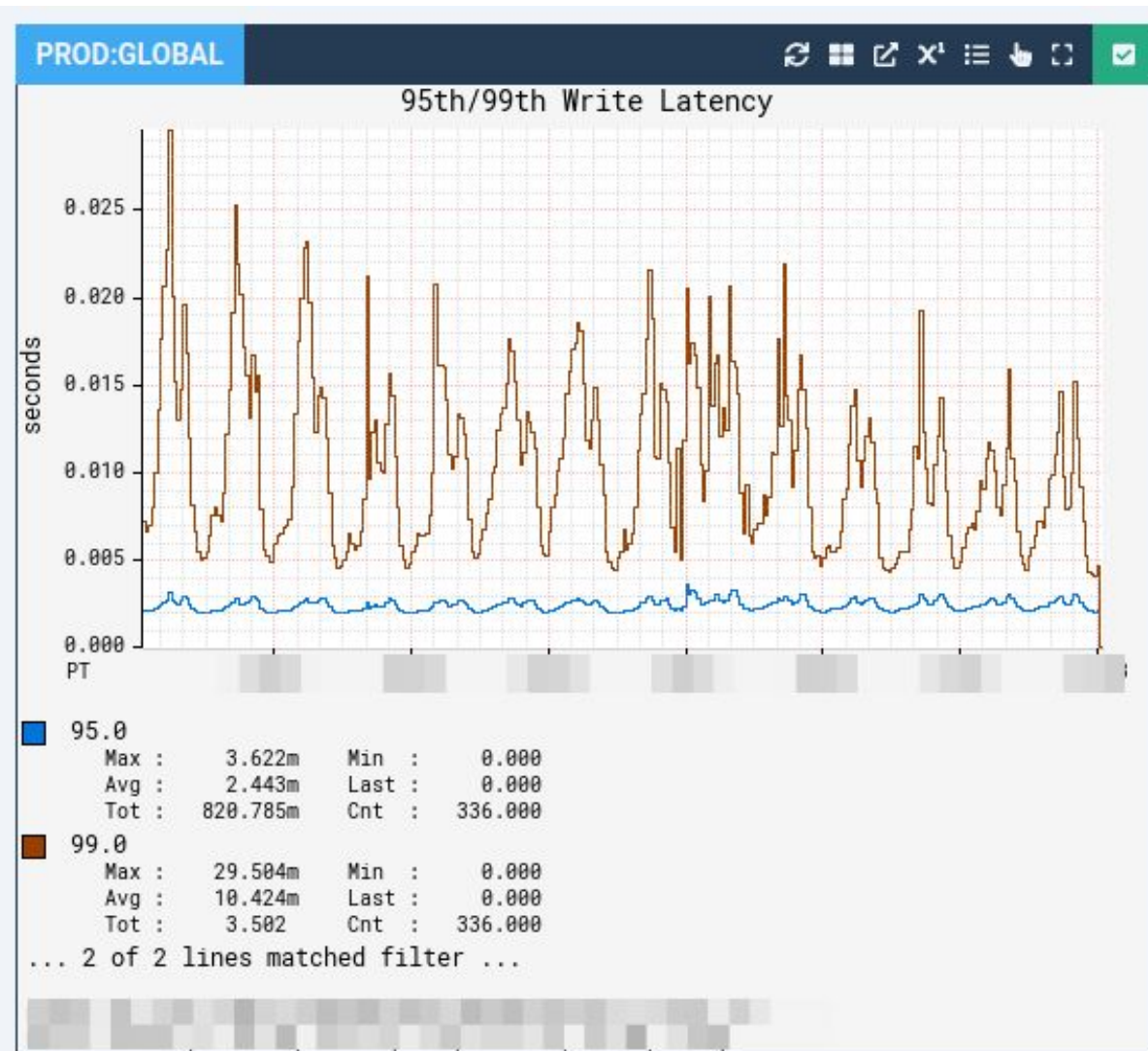
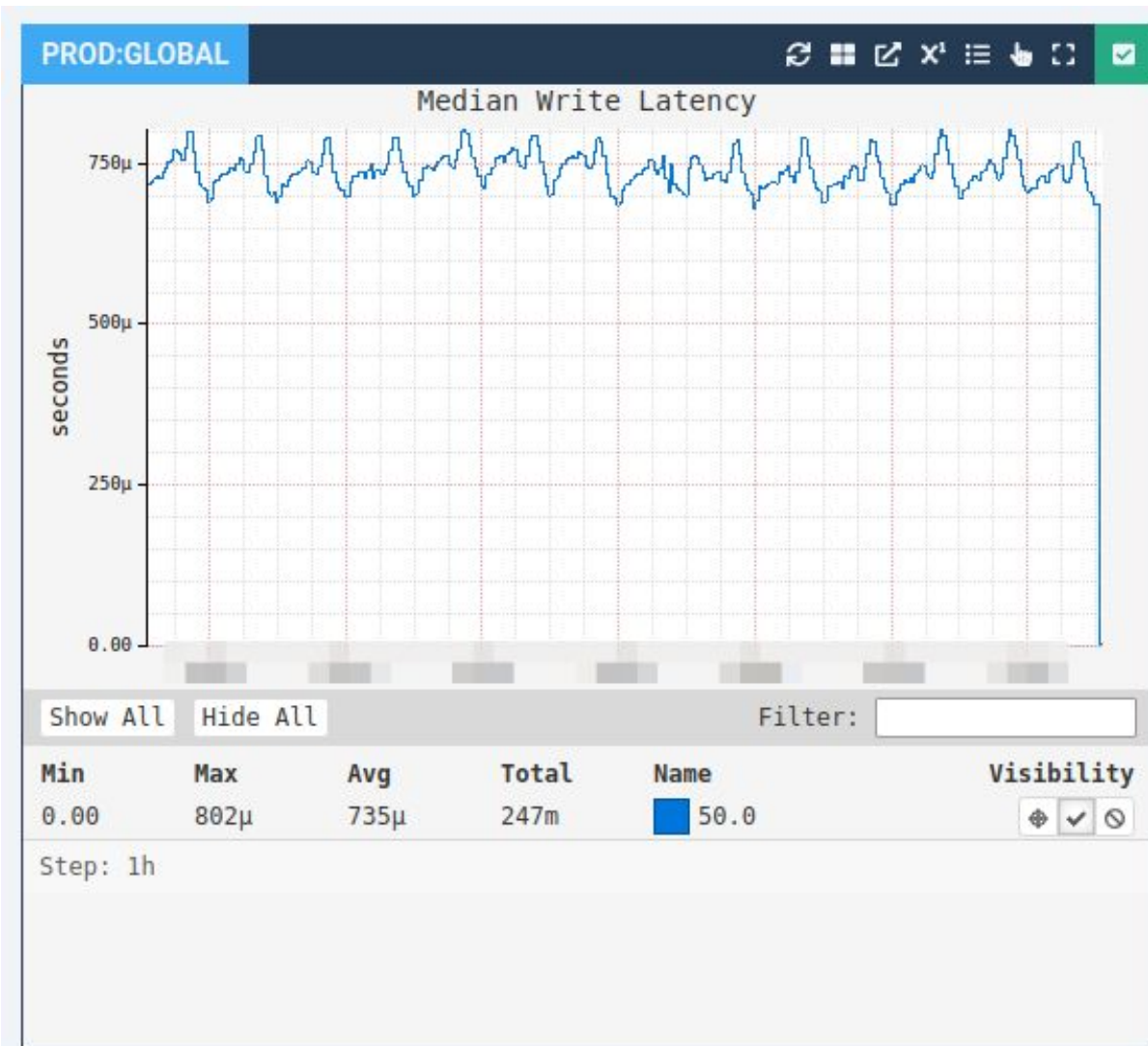# Scale?



**Read Rate**

Peak Traffic is 6 Million Reads per Second

# Scale?

# Scale?

# Conclusions

1. Stay in Zone, failover when loaded

2. LO is easier to load balance for than LQ because we control the entire flow (snitch impacts LQ)

3. We can simulate slow coordinators, and protect against them.

WLLLB is widely deployed at Netflix handling over 10M QPS

# Q/A