AWS
re:Invent

DECEMBER 2 – 6, 2024 | LAS VEGAS, NV

NFX301

# How Netflix handles sudden load spikes in the cloud

**Rob Gulewich**

(he/him)
Principal Software Engineer
Netflix

**Ryan Schroeder**

(he/him)
Staff Software Engineer
Netflix

**Joseph Lynch**

(he/him)
Principal Software Engineer
Netflix

**Manju Prasad**

(she/her)
Sr. Solutions Architect
Amazon

**Rob Gulewich**

Principal Software Engineer
Netflix Platform

**Ryan Schroeder**

Staff Software Engineer
Netflix Reliability

**Joseph Lynch**
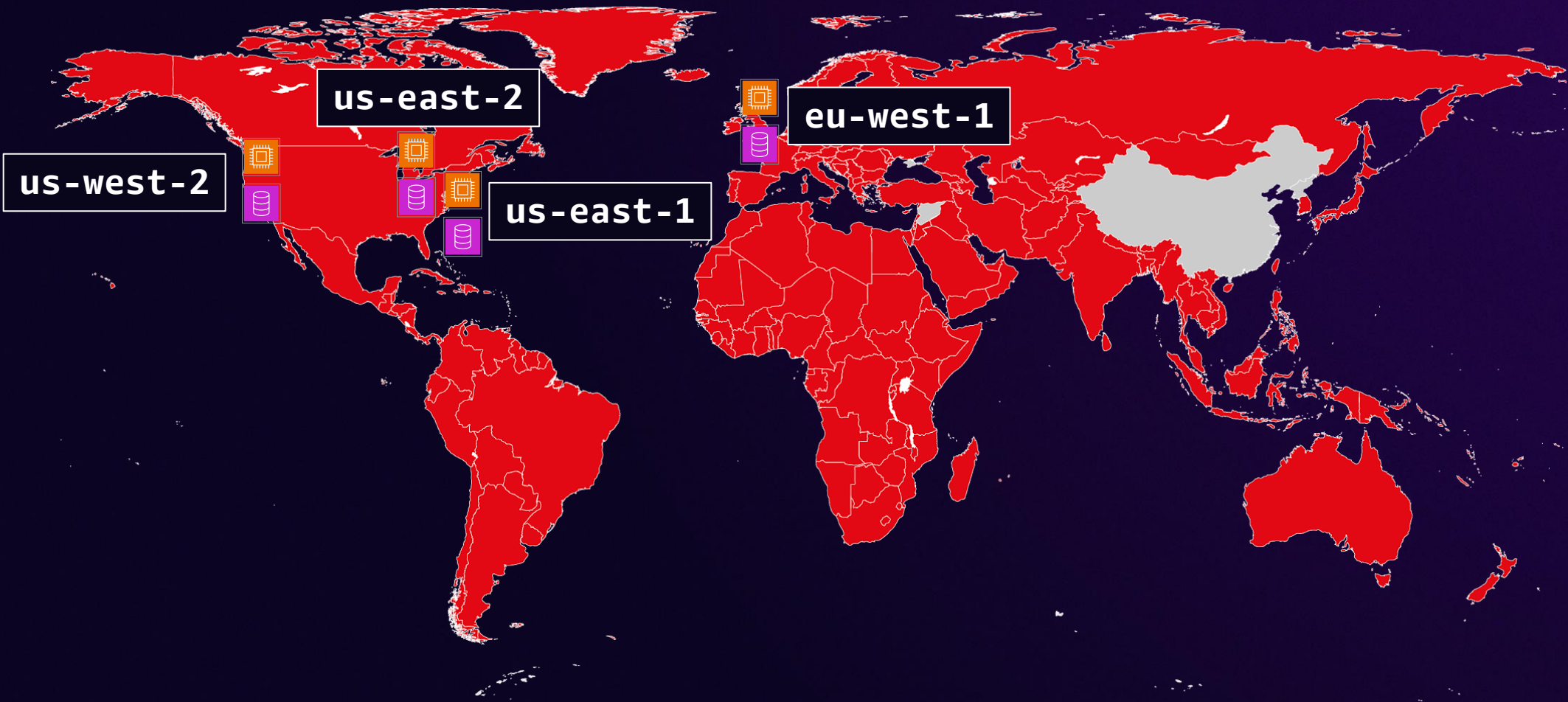
Principal Software Engineer
Netflix Data Platform

# Agenda

**01** Problem: Load spikes

**02** Solution: Predict and plan

**03** Solution: React quickly

**04** Solution: Stay available

**05** Experiment: Test resilience
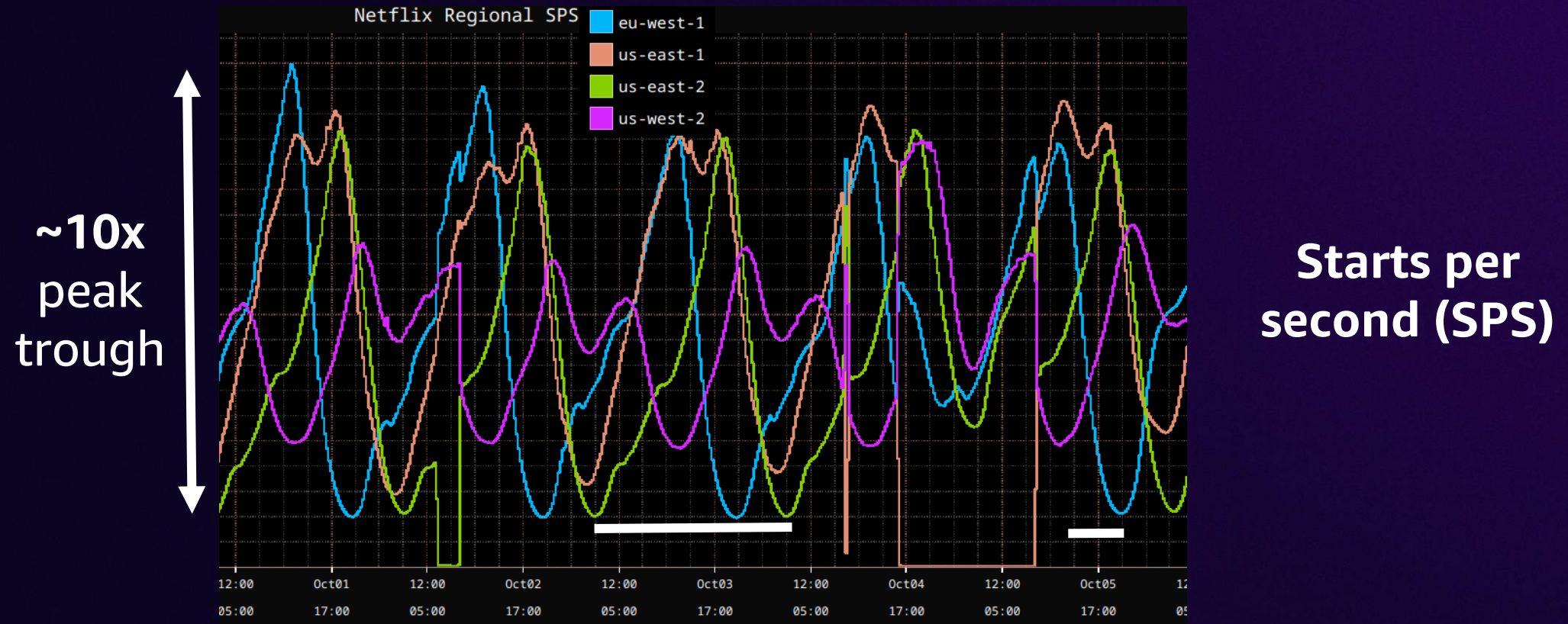
**06** Conclusions and wrap up

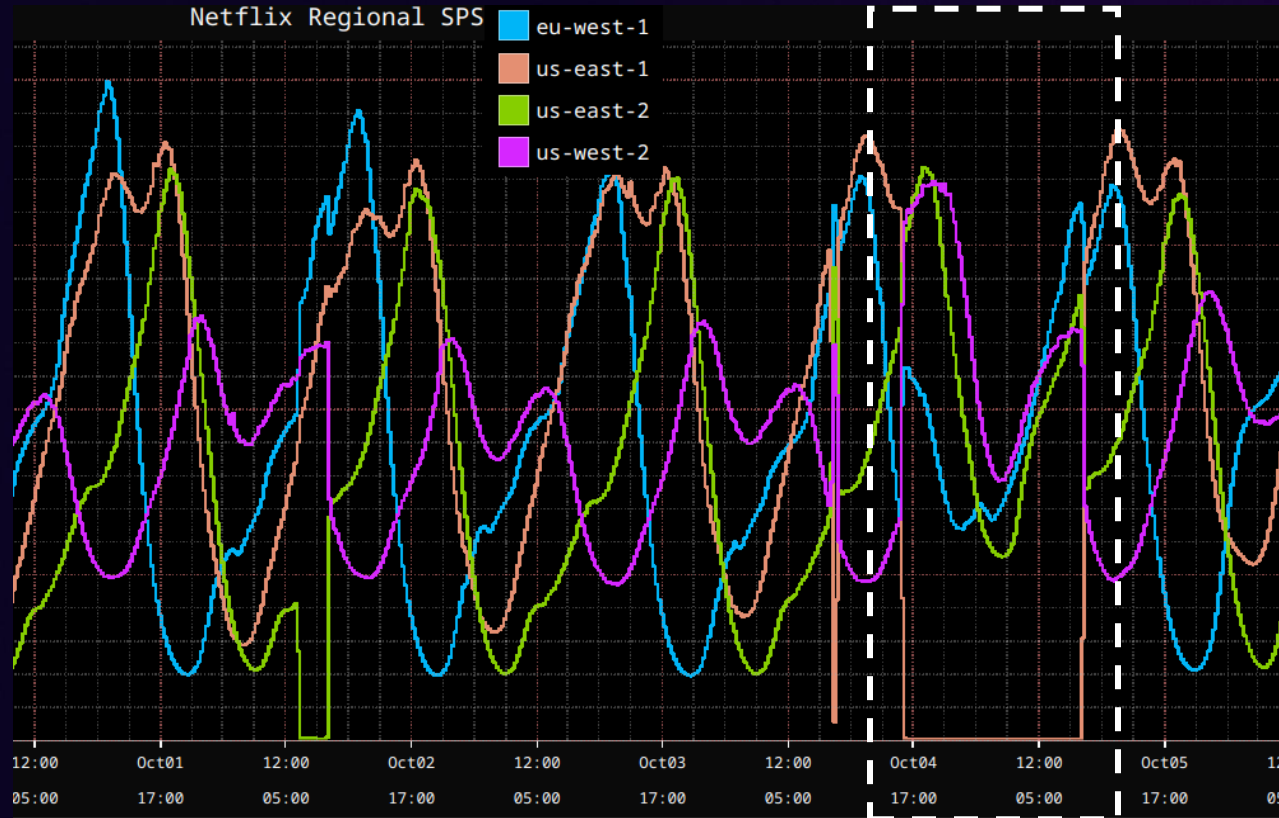# 01: Problem
Load spikes at Netflix

# Netflix cloud topology

# Gradual traffic increases are the norm



~10x peak trough

Starts per second (SPS)
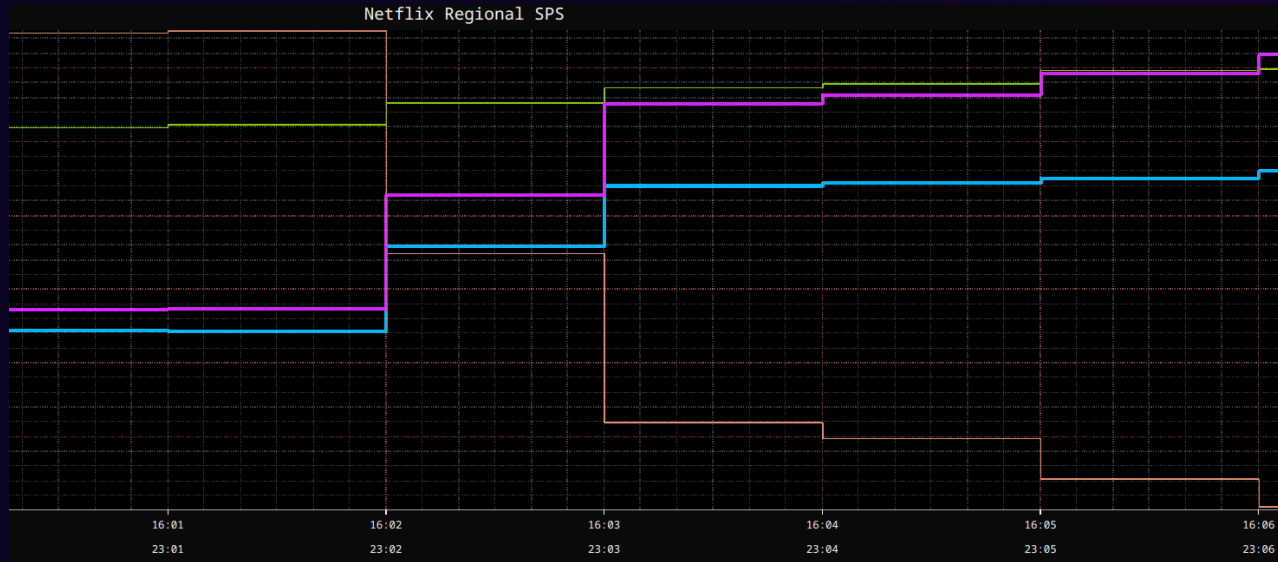
24-hour periodicity

Traffic phase shifts

# Load spikes are common



**Failovers due to:**
* Regular practice
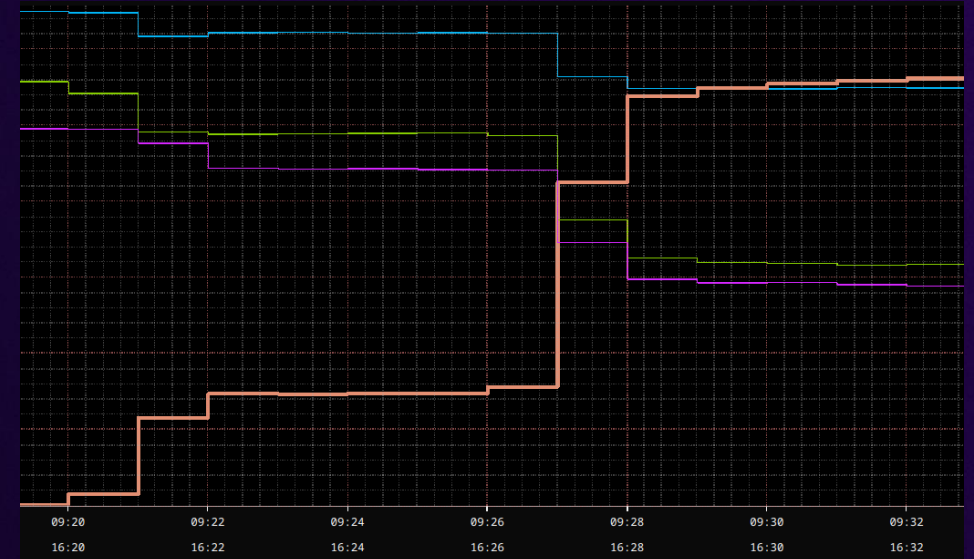* Bad software deployment
* Regional impairment

# Fast failover is a load spike



Netflix Regional SPS

**On Evacuation**
Up to *2x traffic to saviors (variable) in 1 minute*
Long tail traffic takes ~5 minutes
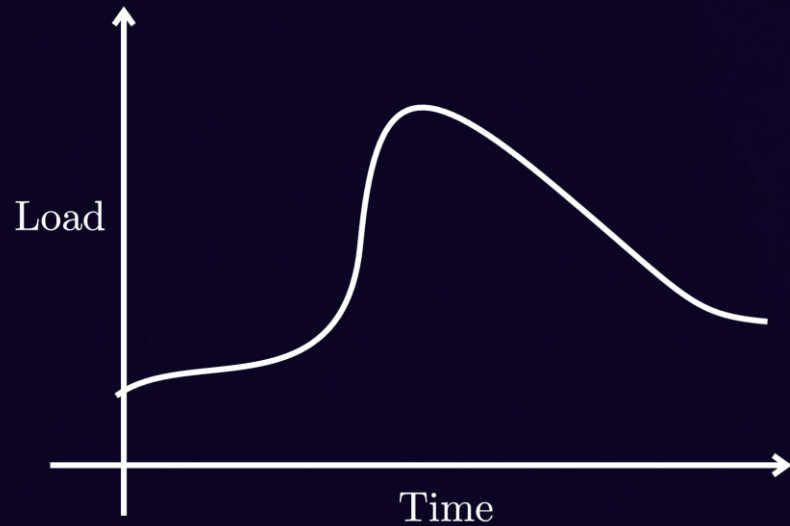Intelligent steering to minimize spike

**On Restore**
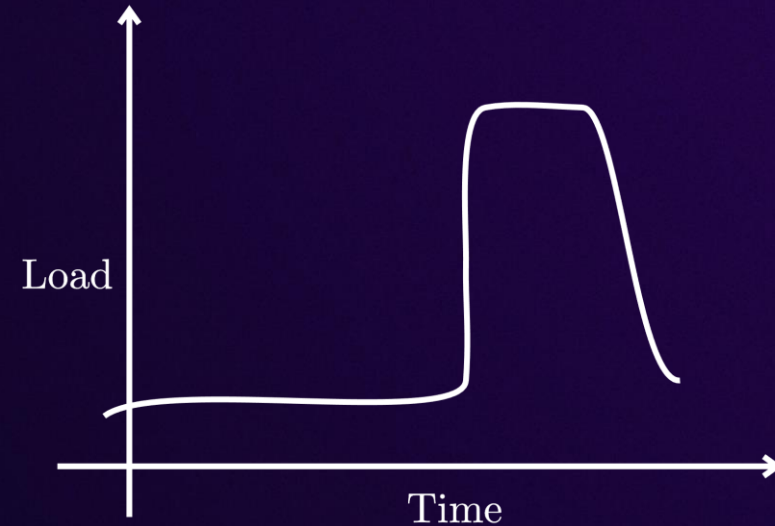Up to *100x traffic to evacuated*

Reintroduce traffic over 5-10 minutes

# Sources of unexpected traffic surges



**Long surges**

Title launches
External events (soccer matches, other sites down)



**Short spikes**

**Retry storms**
**Device bugs**

# Thousands of microservices – Complex downstream call graph

**Microservice-Specific Regional Demand**

Because of service decomposition, we understood that using a proxy demand metric like SPS wasn't tenable and we needed to transition to microservice-specific demand. Unfortunately, due to the diversity of services, a mix of Java (Governator/Springboot with Ribbon/gRPC, etc.) and Node (NodeQuark), there wasn't a single demand metric we could rely on to cover all use cases. To address this, we built a system that allows us to associate each microservice with metrics that represent their demand.
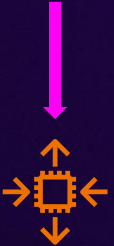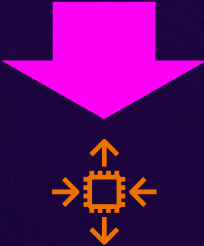
Blog post:

Regional
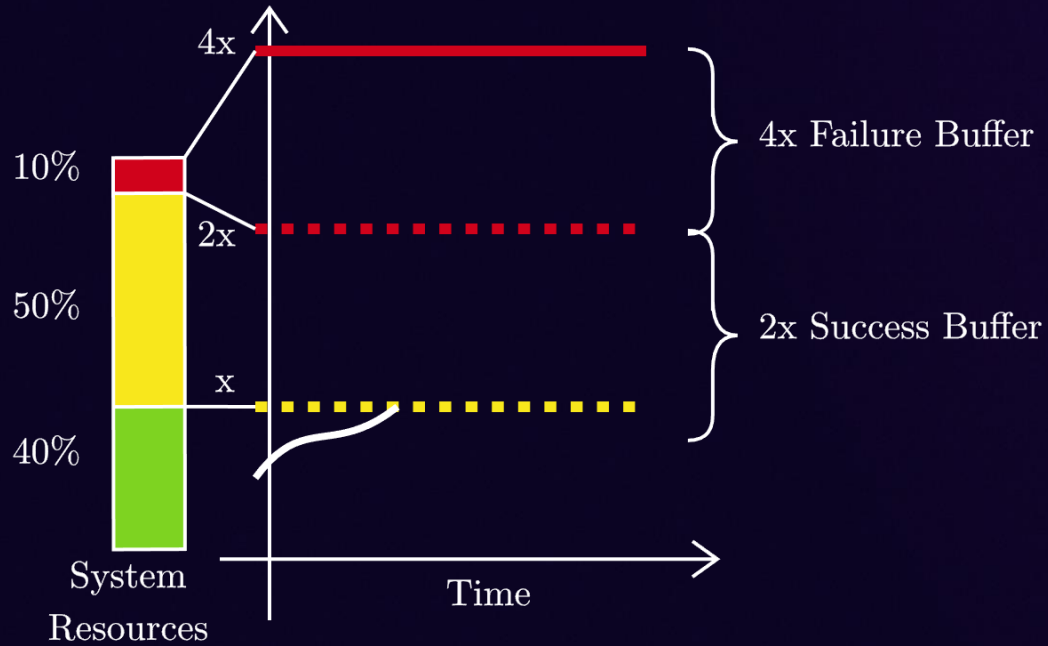2x Spike

Svc A
2x
Tier=0

Svc B
1.2x
Tier=1

Svc C
4x
Tier=2

# Thousands of microservices – Different headroom



Normal System Load with Buffer

System Load Under Load Spike

Every service operates with two key **Buffers**

**Success buffer** Headroom before errors (bad)
**Failure buffer** Headroom before congestive collapse (very bad)

# Our business is evolving

Buffer Recovering After Load Spike



**Changing business needs:**
- More frequent big title launches
- More global launches

**Goals:**
- Reduce time to recover
- Use regional failover less as the primary remediation
- Build resiliency assuming load spikes are the norm

# 02: Capacity plan
## Load is often predictable

# Scale on a schedule

- If we know when traffic is going to arrive, prescale services beforehand to match the predicted load

- Autoscaling is designed for reactive scaling of individual services



Prescale    Spike    Downscale

# Prescaling

- Use the failover system to scale up the entire streaming fleet
- Maps regional SPS to RPS per instance and calculates new min. sizes



Autoscaling Group:
- Blue: min
- Brown: desired
- Green: instances up

**Increase mins to match expected spike**

**Decrease mins to normal levels**

# Shape on a schedule

- Some title launches are centered in a specific geography
- For large launches, we can proactively steer users to other regions to balance global capacity usage

# 03: Scale out of trouble
Predictions are often wrong

# Autoscaling during steady state

# Autoscaling during steady state



Traffic vs Capacity

Smooth Increase

Smooth Decrease

Axis 0
■ Total Traffic
Axis 1
■ Total Capacity

Frame: 2d, End: 2024-10-15T05:27-07:00[US/Pacific], Step: 3m
Fetch: 358ms (L: 8.5M, 4.0k, 2.0; D: 507.0M, 3.9M, 1.9M)

# Autoscaling during load spikes



## Load Spike vs Autoscaling

**Axis 0**

■ Requests Per Second

| | | | |
|---|---|---|---|
| Max : | 137.646k | Min : | 816.533 |
| Avg : | 59.312k | Last : | 59.924k |
| Tot : | 2.372M | Cnt : | 40.000 |

**Axis 1**

■ Autoscaling Group Desired Size

| | | | |
|---|---|---|---|
| Max : | 57.000 | Min : | 5.000 |
| Avg : | 24.375 | Last : | 57.000 |
| Tot : | 975.000 | Cnt : | 40.000 |

Frame: 40m, End: 2024-10-03T17:01-07:00[US/Pacific], Step: 1m
Fetch: 18ms (L: 8.2k, 1.6k, 2.0; D: 490.5k, 64.2k, 80.0k)

# Autoscaling during load spikes

# Components of time-to-recovery



Traffic spikes at $T_L$ causing Utilization to increase. The Cluster Size increases only after delays for Detection and Control Plane. After a delay for O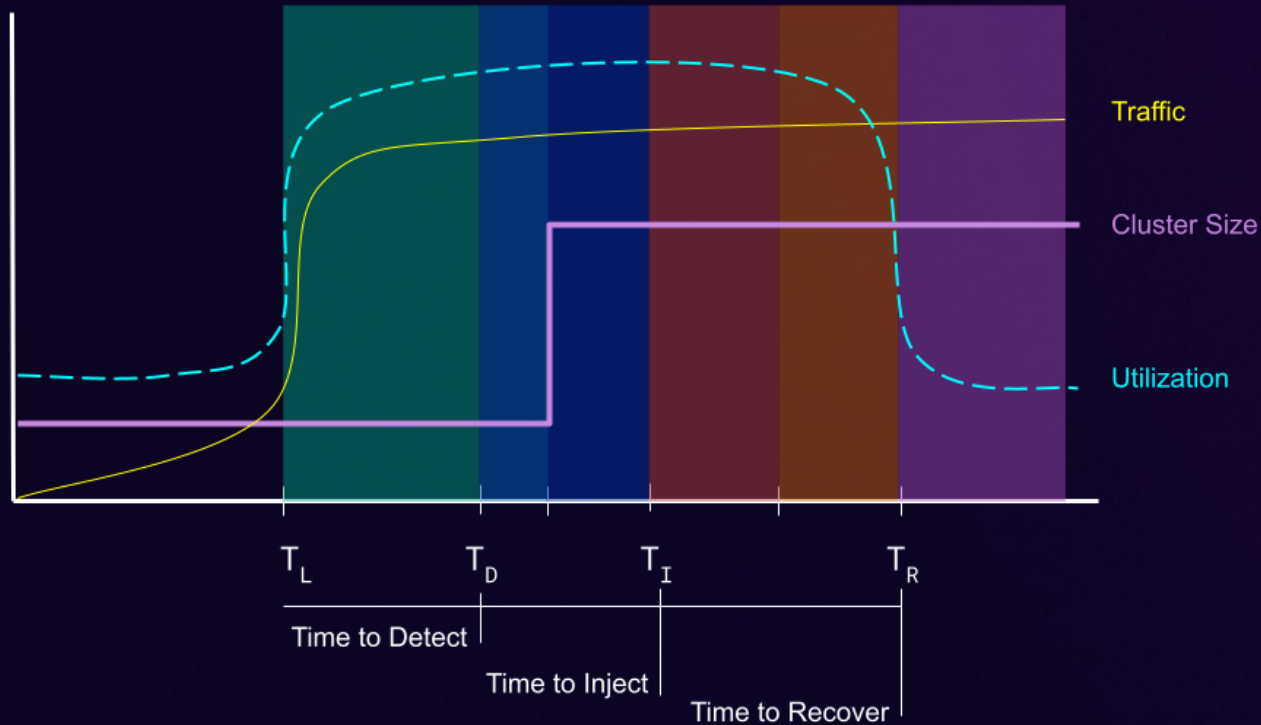S Startup, we reach the point usable capacity is injected $T_I$. Utilization remains high until Application Startup and Load Balancing delays allow new capacity to take traffic - then we Recover at $T_R$.

| Stage | Description |
|---|---|
| Detection | Scaling alarm triggers |
| Control plane | Hardware online |
| System startup | Kernel and base systemd units started |
| Application startup | Microservice started |
| Traffic | Traffic arrives |

# Experimentation setup

- Synthetic load generation
- Baseline vs. experiment comparison
- Variations of scaling policy configurations

# Time-to-recovery dominating factors



Breakdown of Startup Latency

Detection > App Startup > System Startup > Hardware Startup

# Time-to-recovery dominating factors
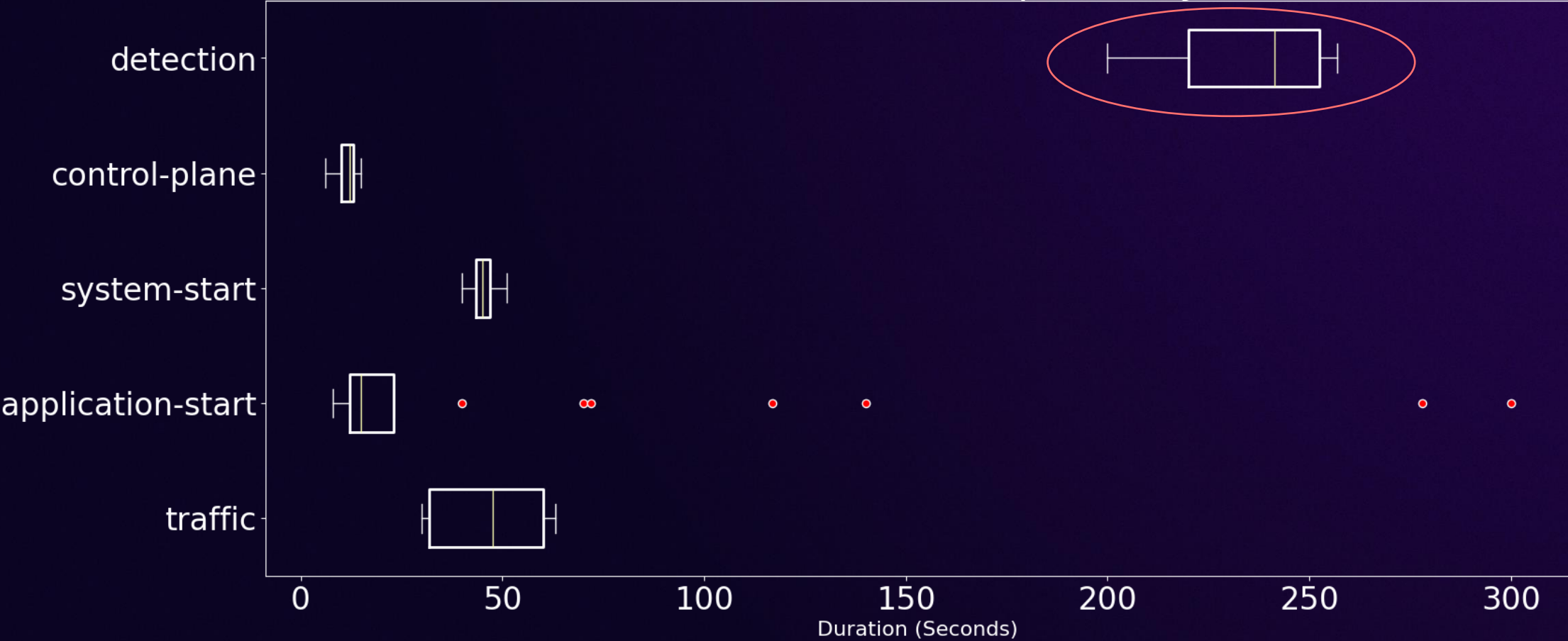


Breakdown of Startup Latency

**Detection** > App Startup > System Startup > Hardware Startup

# Detection – Scaling on RPS

CPU target tracking is nice for gradual changes, but doesn't provide enough information for 10x spikes

- Typical CPU target is around 50% utilization

- At 2x RPS, CPU is 100%

- At 10x RPS, CPU is also 100%

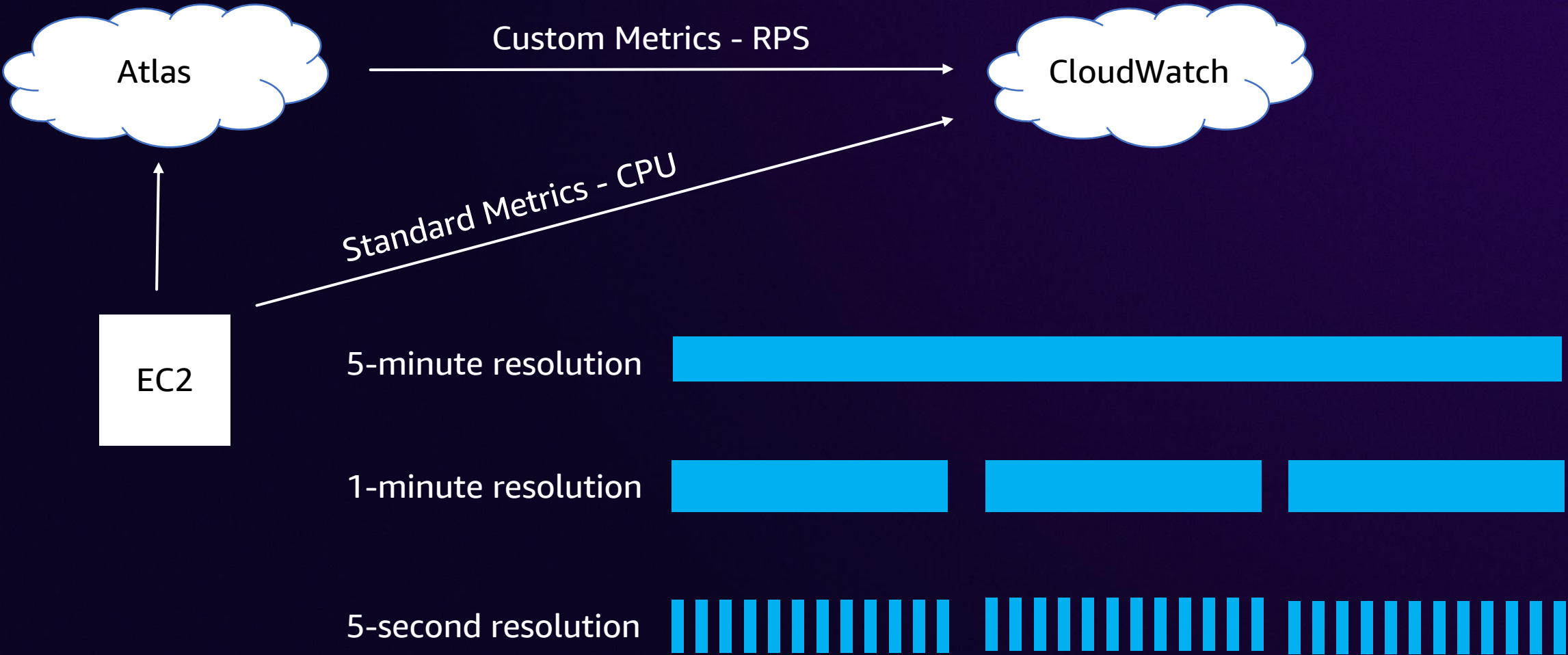During load shedding, CPU does not reflect actual workload

CPU:   0% – 100%

RPS:   0 – Infinity

# Detection – Scaling on RPS

Add RPS "hammer" policy – one shot to success

- Bad – 2x scales
- Good – exactly what you need
- Bad – scale way too much

# Detection – Higher resolution metrics

Atlas

CloudWatch

Custom Metrics - RPS

Standard Metrics - CPU

EC2

5-minute resolution

1-minute resolution

5-second resolution

# Time-to-recovery dominating factors



Breakdown of Startup Latency

**Detection** > App startup > System startup > Hardware startup

# Time-to-recovery dominating factors



Breakdown of Startup Latency

3x Improvement

Detection > **App startup** > System startup > Hardware startup

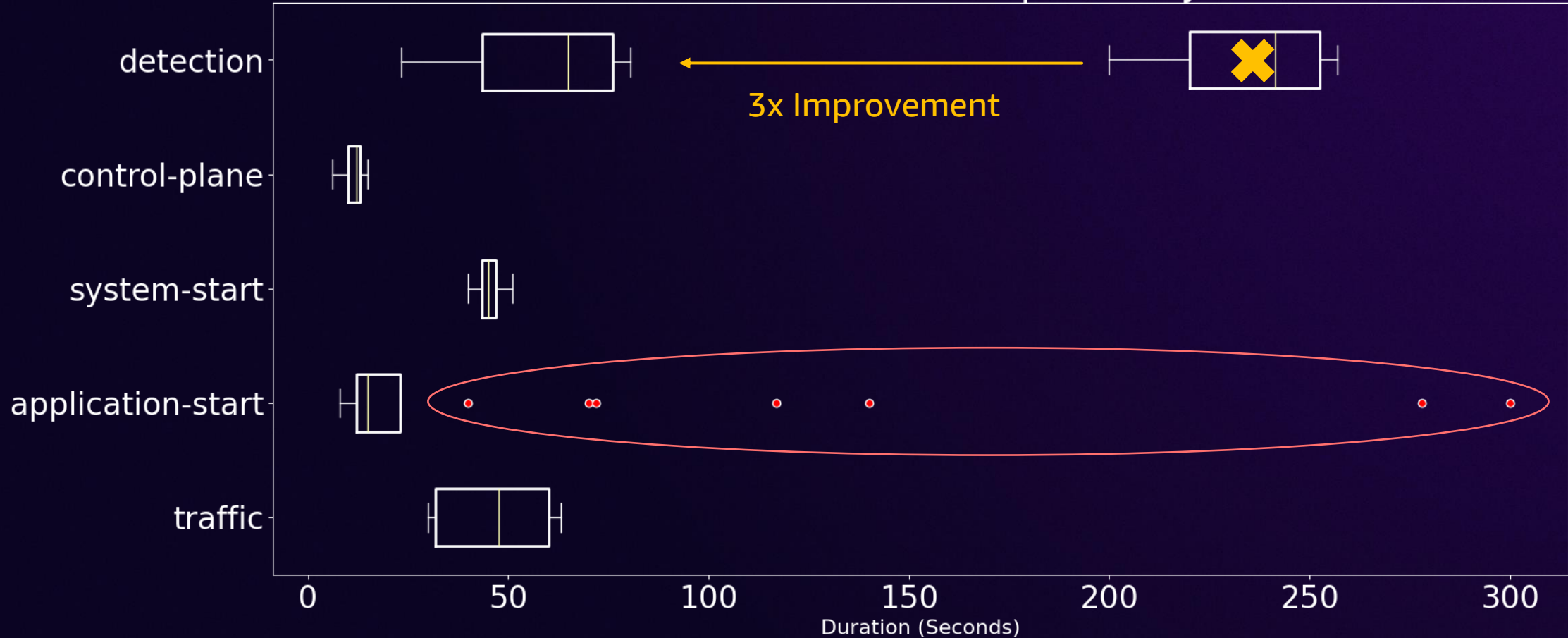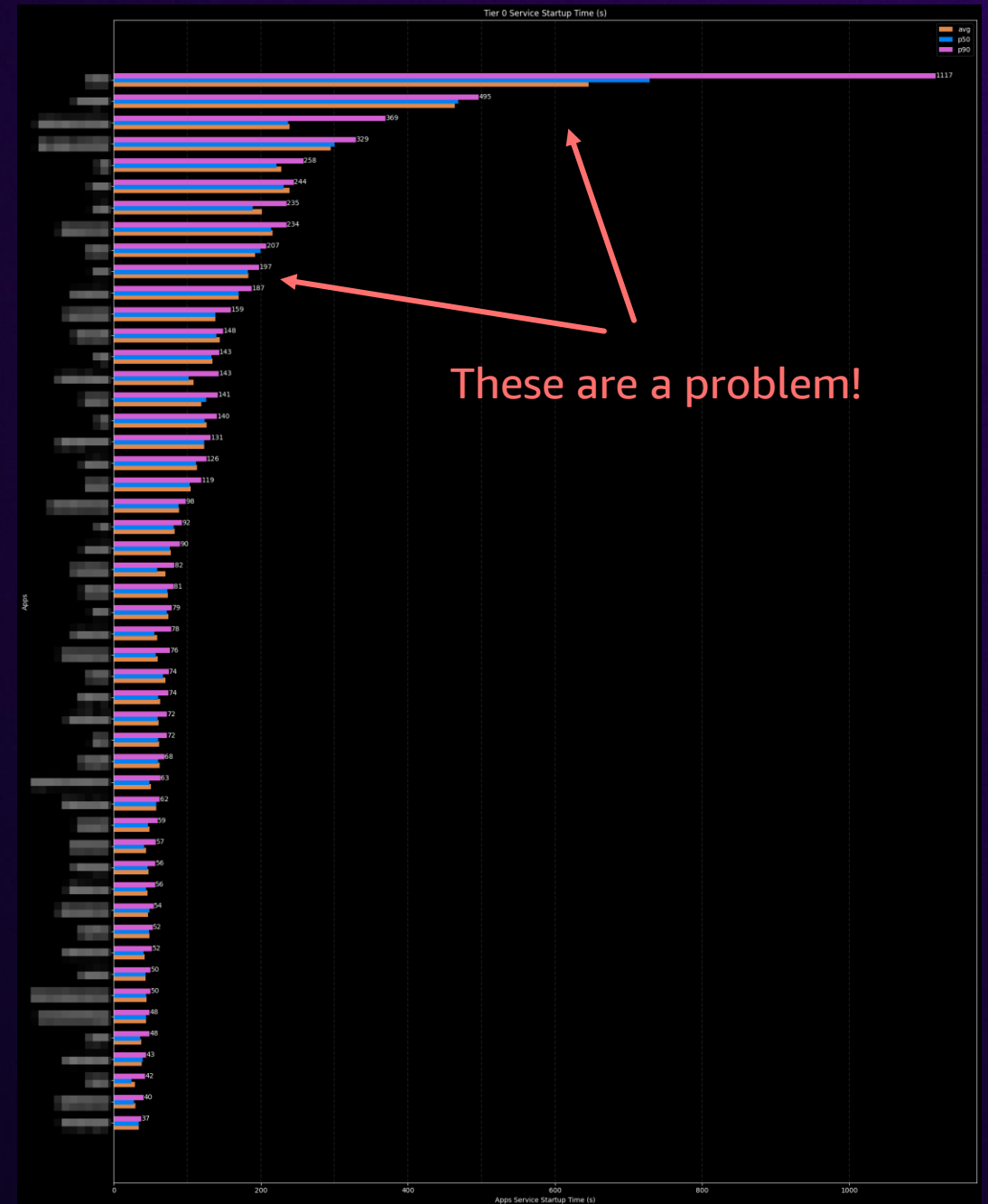# App startup has a long tail

Long tail of startup delay. Vast majority under a minute.

Worst offender took **p90 of *18 minutes*** to start!



These are a problem!

# Time-to-recovery dominating factors



Breakdown of Startup Latency

3x Improvement

Eliminated the Long Tail

Detection > **App startup** > System startup > Hardware startup

# Time-to-recovery dominating factors

Breakdown of Startup Latency



Detection > App startup > **System startup** > Hardware startup

# Start system faster – `systemd-analyze`
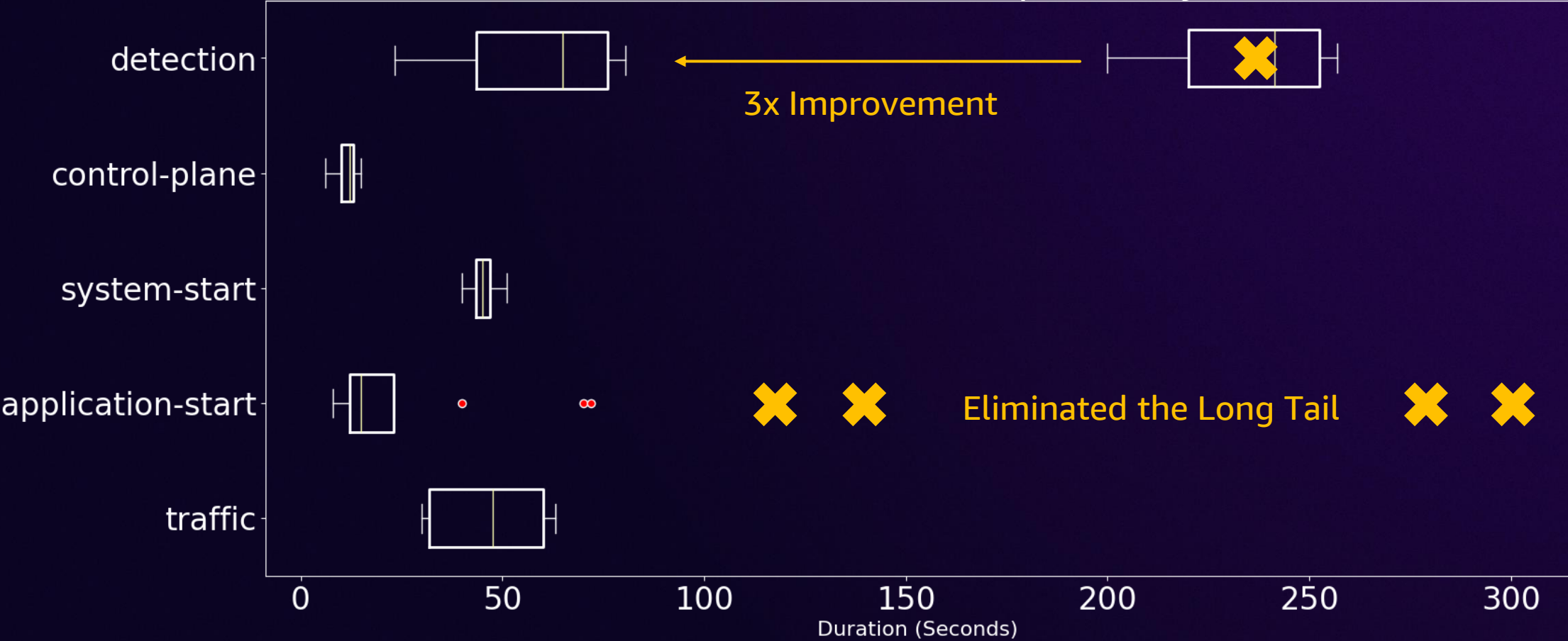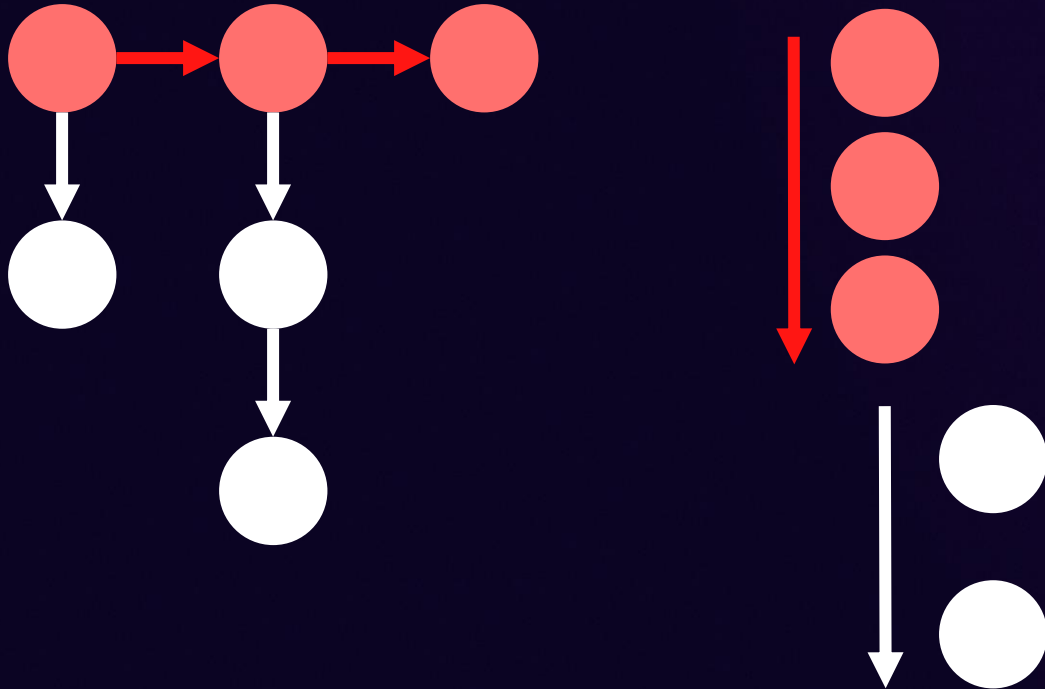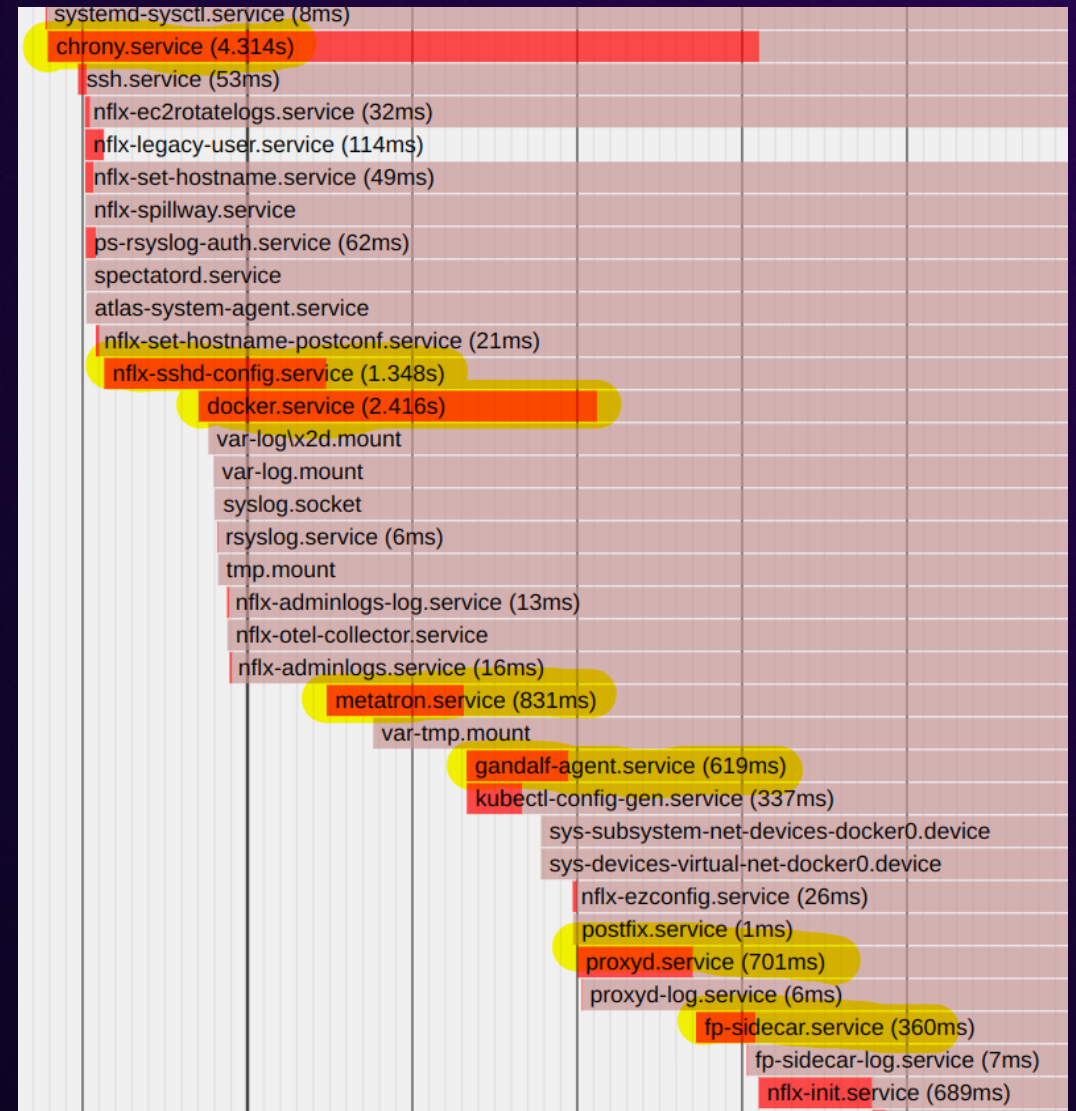
**Sequential –**
*Slow*

**Parallel -**
*Fast*



systemd-sysctl.service (8ms)
chrony.service (4.314s)
ssh.service (53ms)
nflx-ec2rotatelogs.service (32ms)
nflx-legacy-user.service (114ms)
nflx-set-hostname.service (49ms)
nflx-spillway.service
ps-rsyslog-auth.service (62ms)
spectatord.service
atlas-system-agent.service
nflx-set-hostname-postconf.service (21ms)
nflx-sshd-config.service (1.348s)
docker.service (2.416s)
var-log\x2d.mount
var-log.mount
syslog.socket
rsyslog.service (6ms)
tmp.mount
nflx-adminlogs-log.service (13ms)
nflx-otel-collector.service
nflx-adminlogs.service (16ms)
metatron.service (831ms)
var-tmp.mount
gandalf-agent.service (619ms)
kubectl-config-gen.service (337ms)
sys-subsystem-net-devices-docker0.device
sys-devices-virtual-net-docker0.device
nflx-ezconfig.service (26ms)
postfix.service (1ms)
proxyd.service (701ms)
proxyd-log.service (6ms)
fp-sidecar.service (360ms)
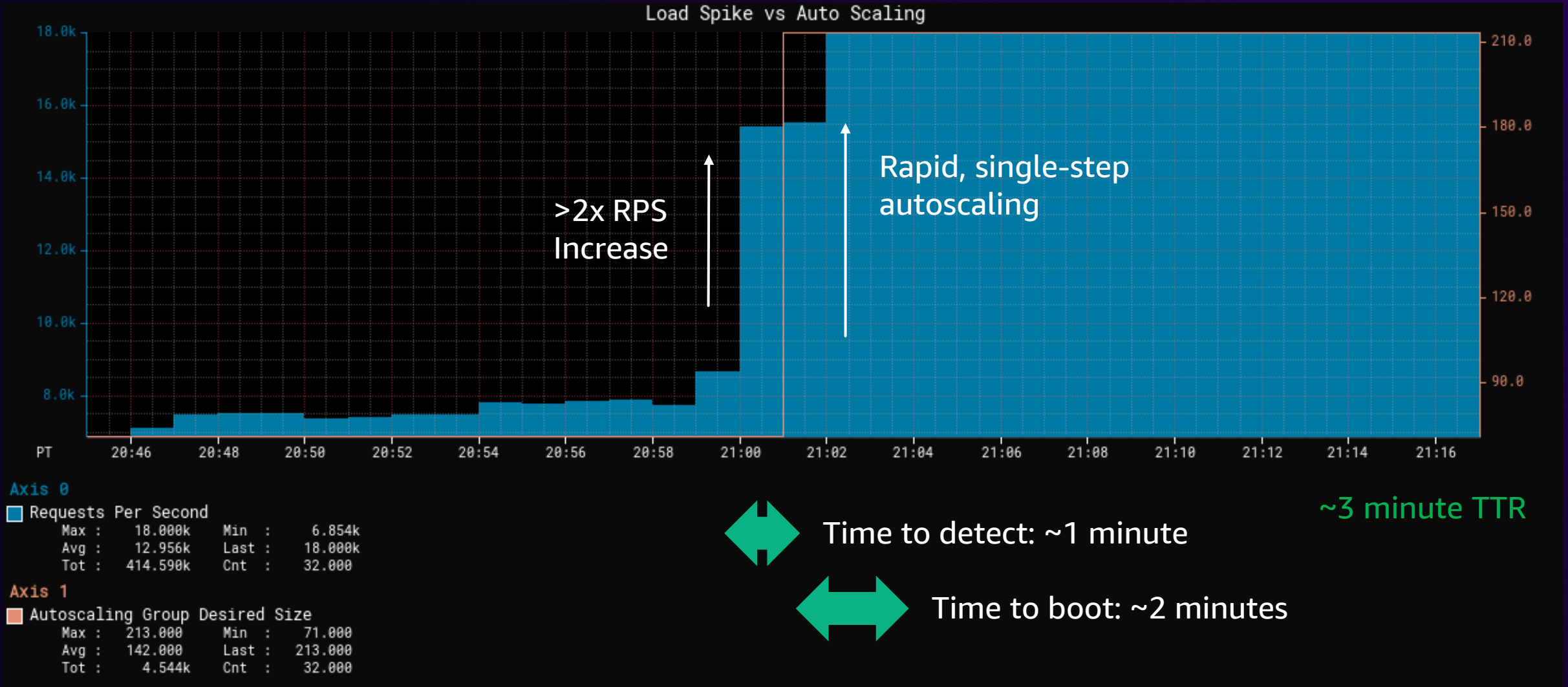fp-sidecar-log.service (7ms)
nflx-init.service (689ms)

# Time-to-recovery dominating factors



Breakdown of Startup Latency

Detection > App startup > **System startup** > Hardware startup

# Results



Load Spike vs Auto Scaling

>2x RPS Increase

Rapid, single-step autoscaling

Axis 0

◼ Requests Per Second
| Max : | 18.000k | Min : | 6.854k |
| Avg : | 12.956k | Last : | 18.000k |
| Tot : | 414.590k | Cnt : | 32.000 |

Axis 1

◼ Autoscaling Group Desired Size
| Max : | 213.000 | Min : | 71.000 |
| Avg : | 142.000 | Last : | 213.000 |
| Tot : | 4.544k | Cnt : | 32.000 |

~3 minute TTR

Time to detect: ~1 minute

Time to boot: ~2 minutes

# 04: Stay available
Techniques to stay up

# Build shared criticality nomenclature

Define some tags, and start tagging

Less debating, more tagging

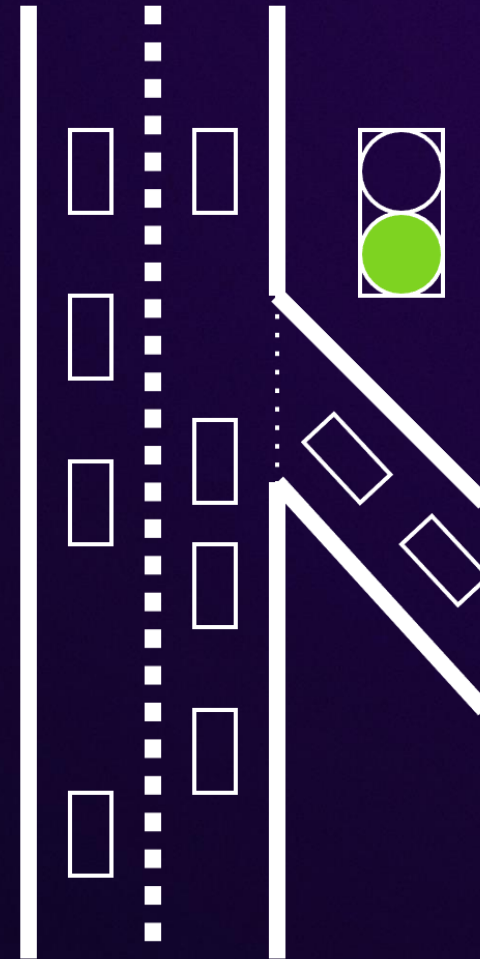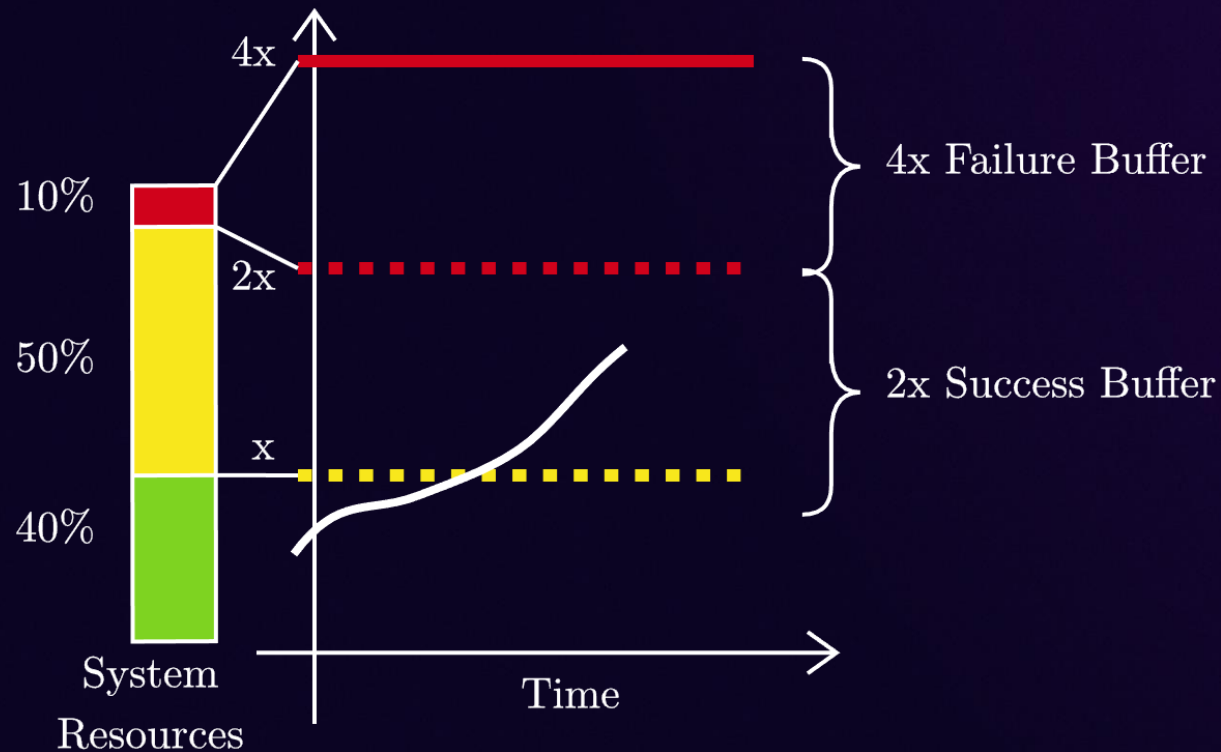| Tag | Values | Reason | Consequence |
|---|---|---|---|
| **Tier** | `int:= {0, 1, 2, 3}` | Spend $$ on what is important | Buffer, testing requirements |
| **Business domain** | `List[str] := {"svod", "gaming", "studio"…}` | Different domains scale differently | Deployment modalities, buffer, … |
| **Lifecycle** | `Str := {"alpha", "beta", "ga", "deprecated", "eol"}` | Do not waste time on deprecated apps | Exclude early/late from requirements |

# Load begins

System Load with Buffer
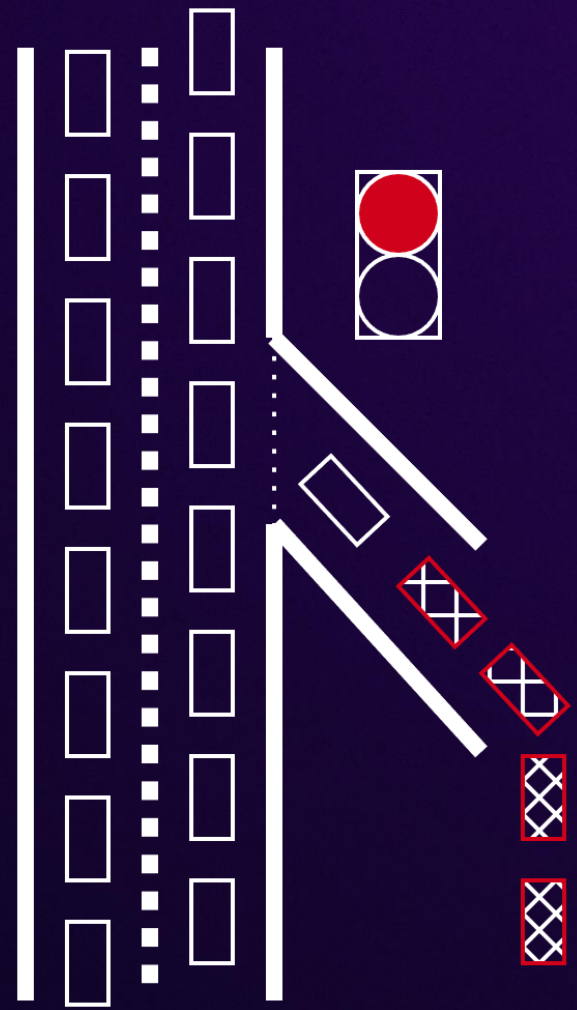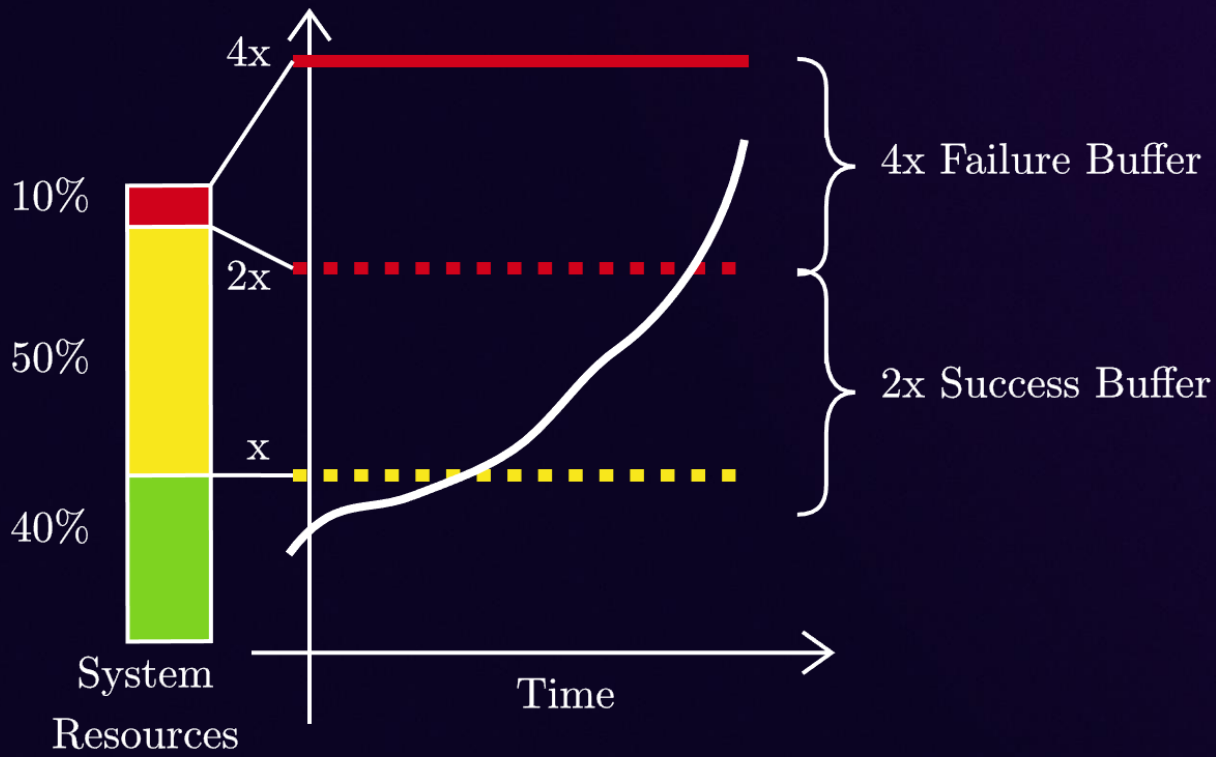


$$T_{transit} = 10m$$

# Load grows

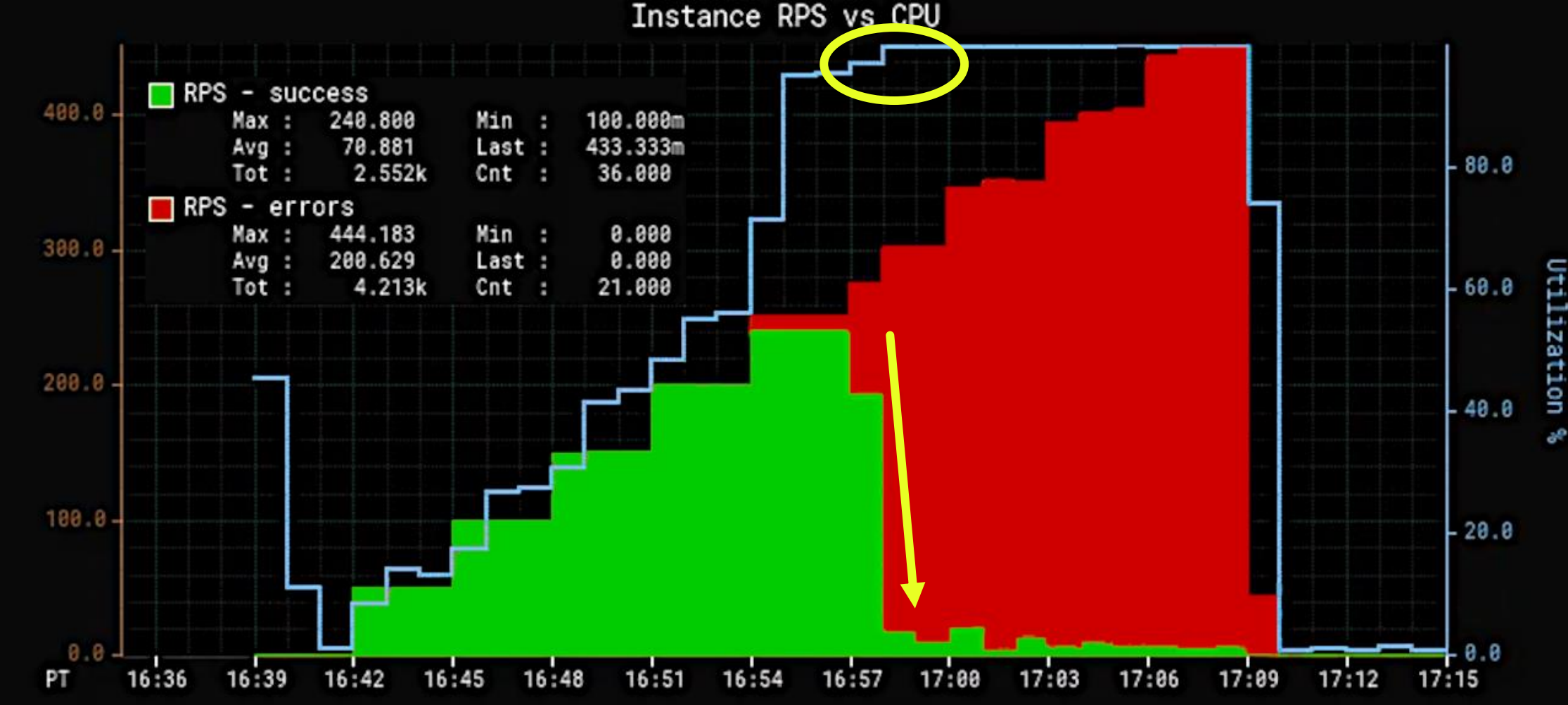## System Load with Buffer



4x

10%

2x

50%

x

40%

System
Resources

Time

4x Failure Buffer

2x Success Buffer



$$T_{transit} = 15m$$

# Load sheds



System Load with Buffer

4x ─── 4x Failure Buffer

2x ····· 2x Success Buffer

x ·····

10%
50%
40%

System Resources

Time

$$T_{transit} = 30m$$

# Congestive failure – very *bad*



System Load with Buffer

4x

10%

2x

50%

x

40%

4x Failure Buffer

2x Success Buffer

System Resources

Time

$$T_{transit} = \infty$$

# Congestive failure

# Stay up – Prioritized shedding

**Success Buffer**

**Prioritized** CPU Shedding


**Failure Buffer**

**Unprioritized** CPU Shedding

Blog post:



**Enhancing Netflix Reliability with Service-Level Prioritized Load Shedding**

Applying Quality of Service techniques at the application level

Netflix Technology Blog · Follow
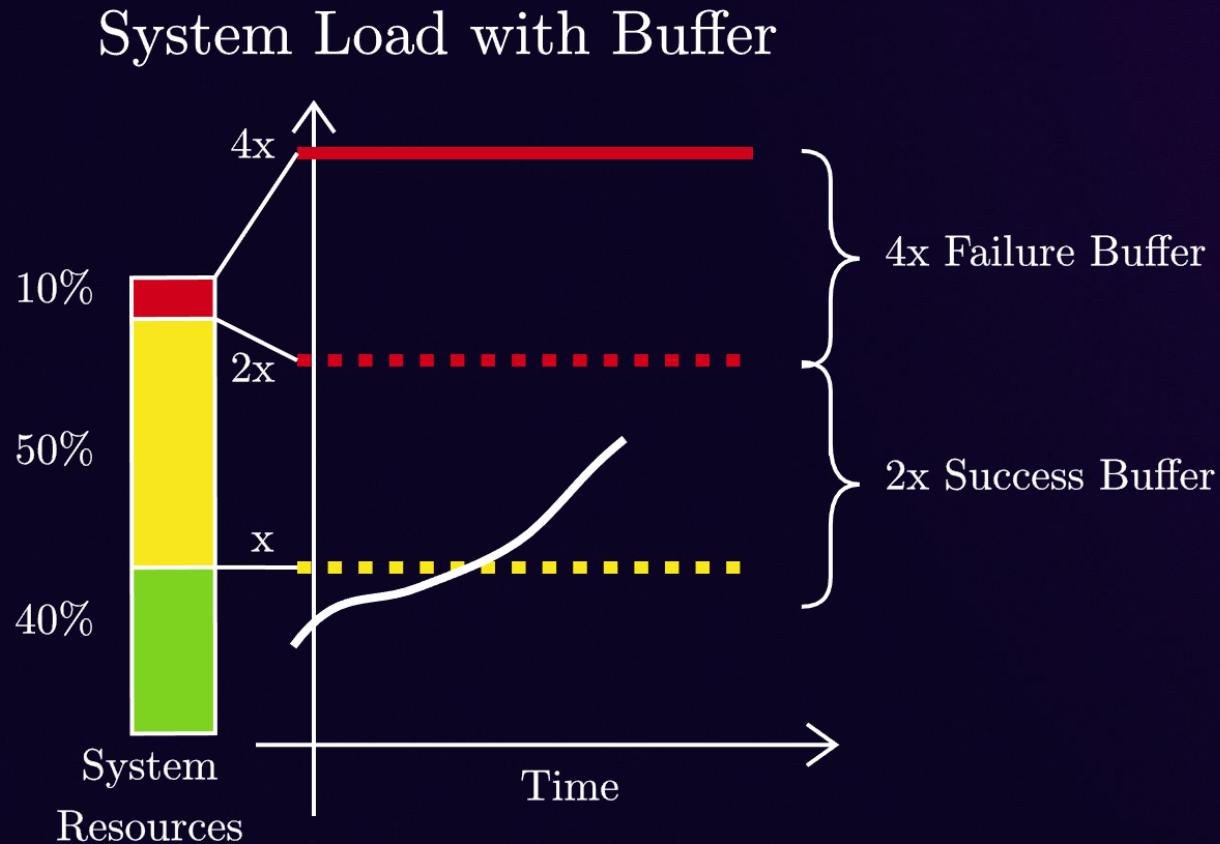Published in Netflix TechBlog · 12 min read · Jun 24, 2024

👏 337    💬 4

Anirudh Mendiratta, Kevin Wang, Joey Lynch, Javier Fernandez-Ivern, Benjamin Fedorka

**Introduction**

In November 2020, we introduced the concept of prioritized load shedding at the API gateway level in our blog post, Keeping Netflix Reliable Using Prioritized Load Shedding. Today, we're excited to dive deeper into how we've extended this strategy to the individual service level, focusing on the video streaming control plane and data plane, to further enhance user experience and system resilience.

# Stay up – Allocate buffers

System Load with Buffer
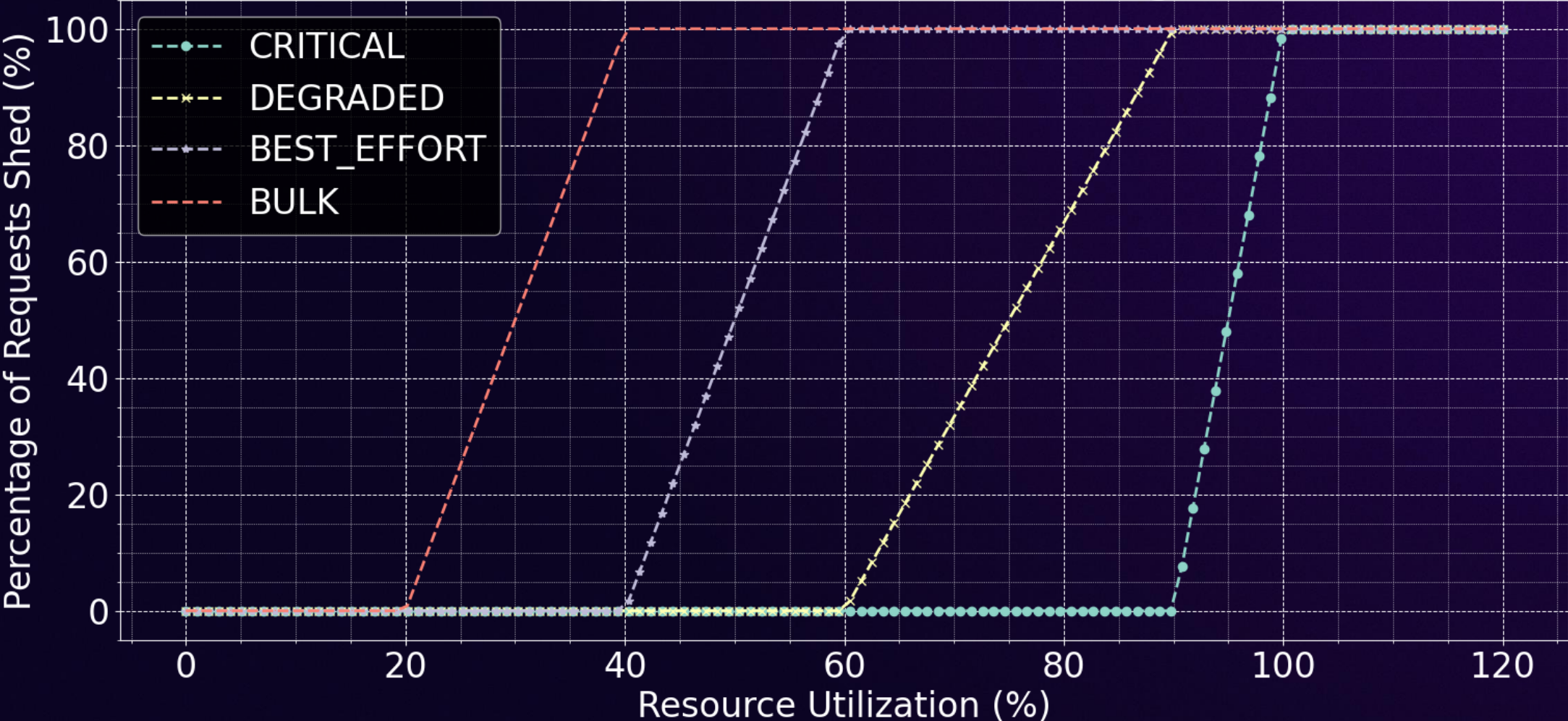


```
resource:
  utilization:
    cpu:
      target: 40
      max: 90
```

$$\text{Buffer}_{\text{success}} \propto T_{\text{startup}}$$
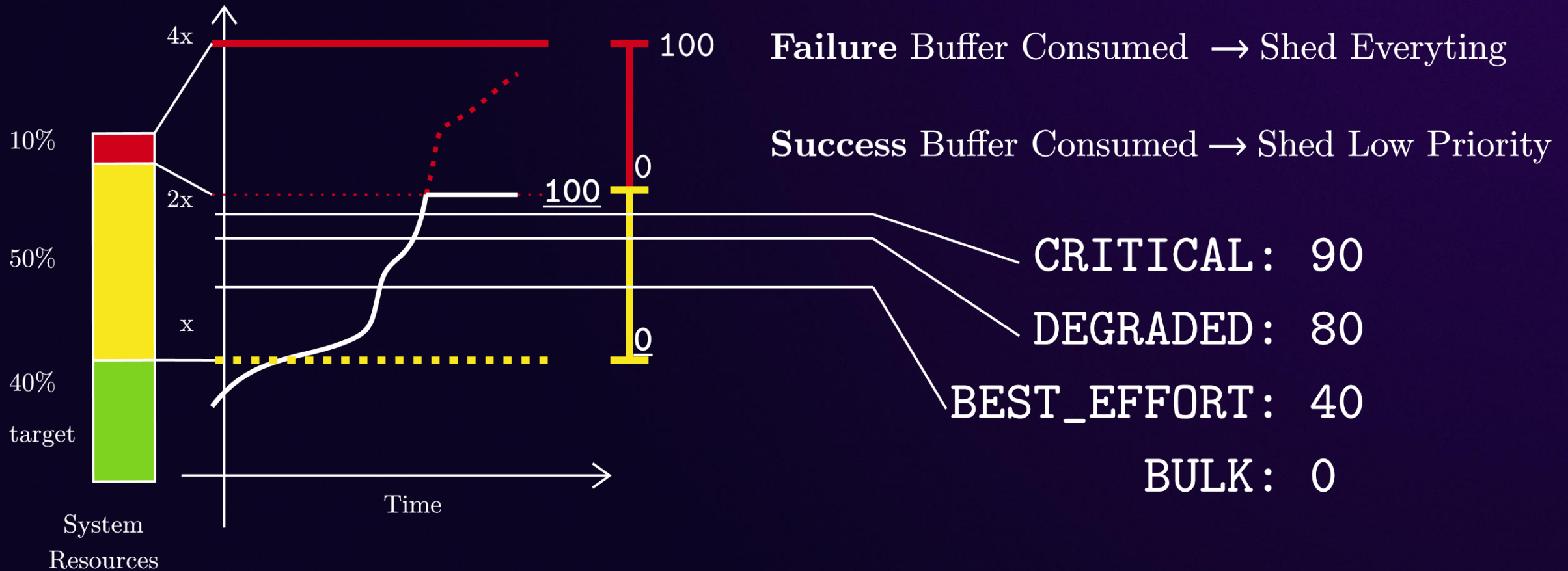
# Stay up – Define priority buckets



Progressive Load Shedding
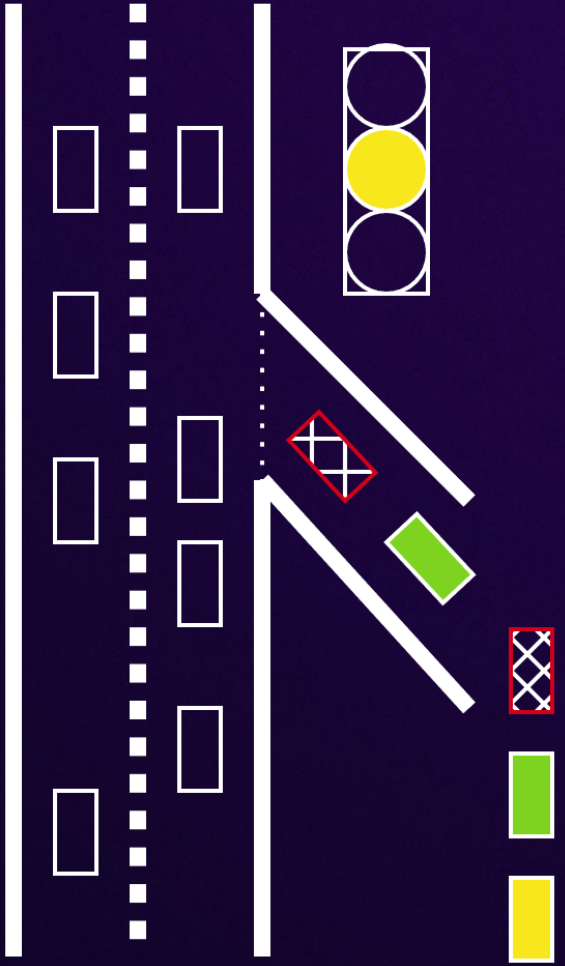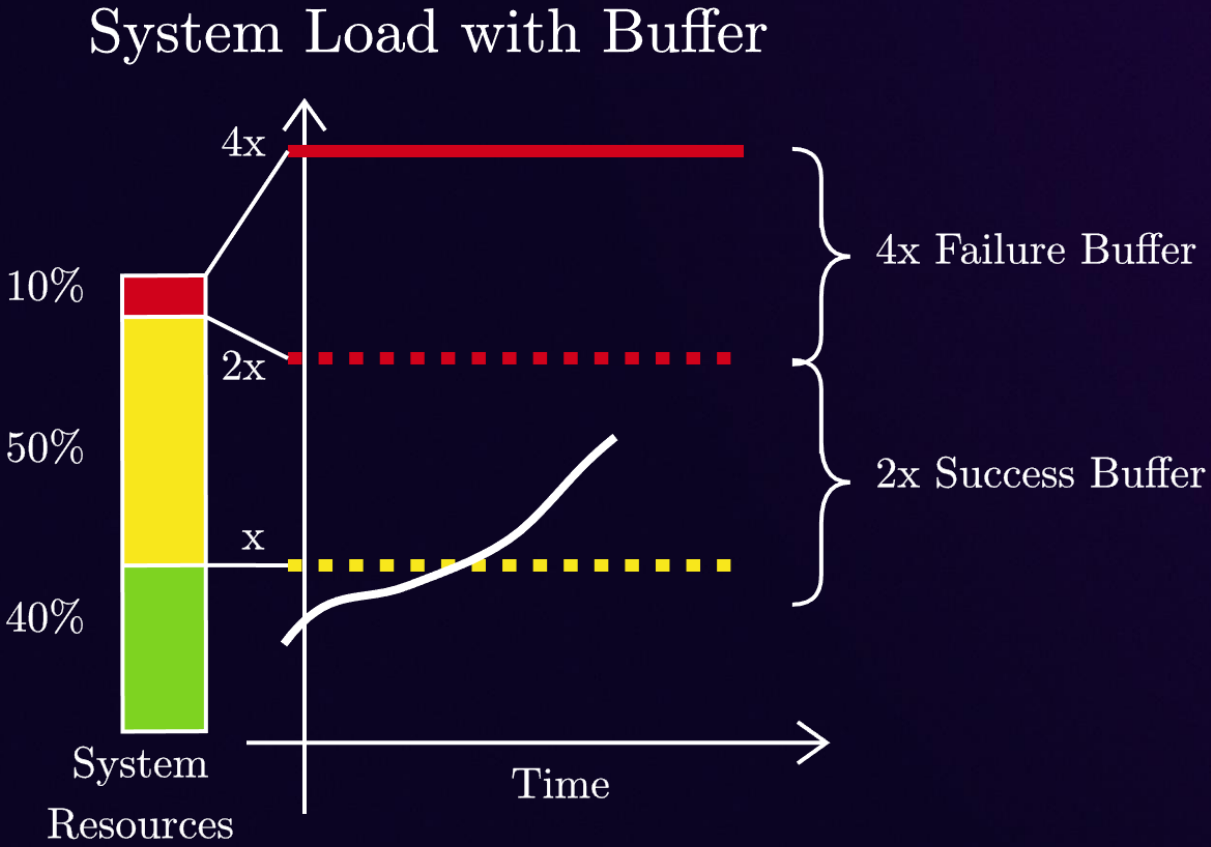
# Stay up – Allocate buffers



System Load Under Load Spike - With Prioritizied Shedding in Success Buffer

**Failure** Buffer Consumed → Shed Everyting

**Success** Buffer Consumed → Shed Low Priority

CRITICAL: 90

DEGRADED: 80

BEST_EFFORT: 40

BULK: 0

# Stay up – Prioritize requests

```java
return context -> {
  Request req = context.getRequest();
  // Prioritize a particular path
  if (req.getPath().startsWith("/critical-play-url")) {
    return PriorityBucket.CRITICAL;
  }

  // Deprioritize background requests
  if (req.getParams().contains("background")) {
    return PriorityBucket.DEGRADED;
  }

  // Take the client device priority
  return getClientPriority(context.getHeaders());
}
```
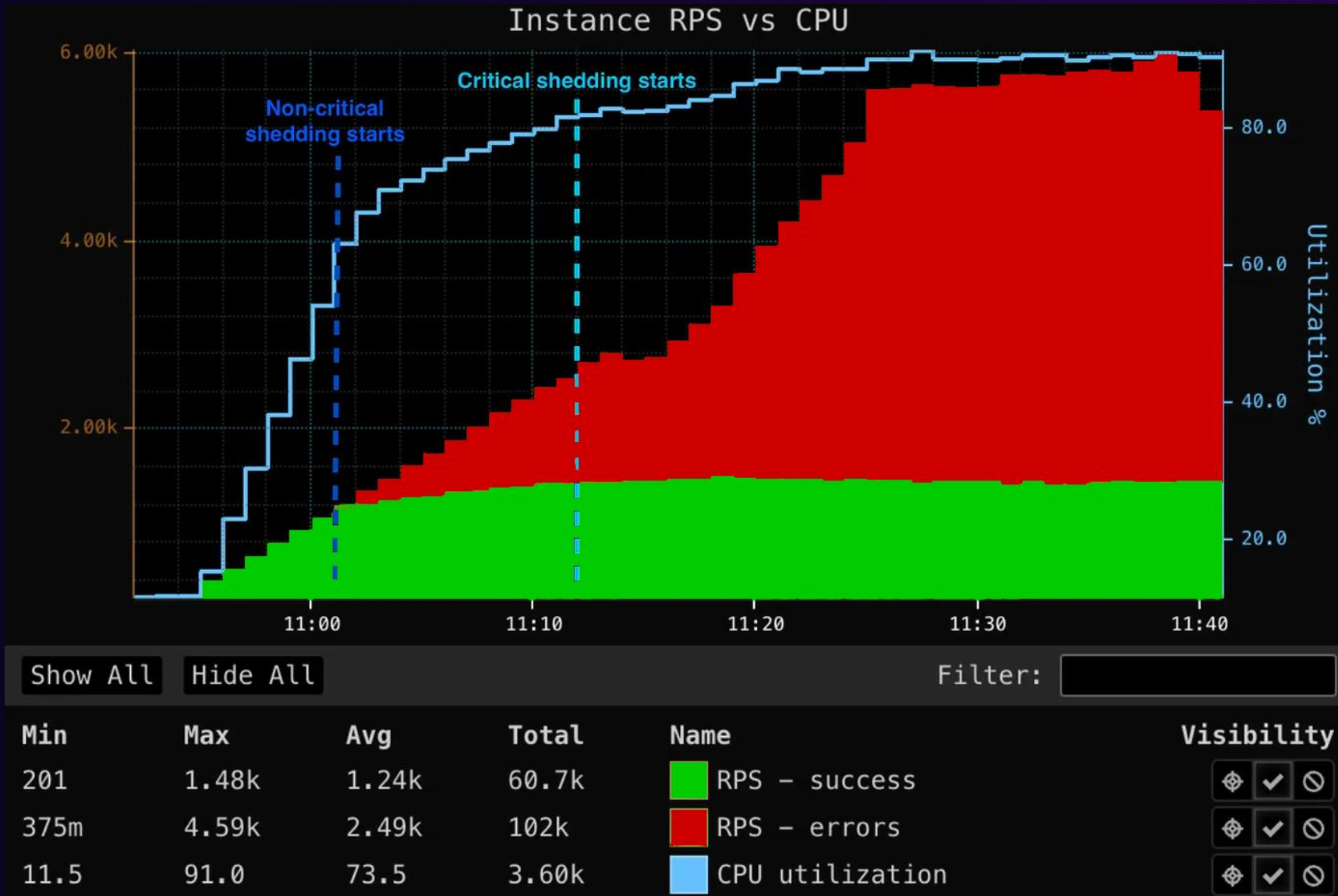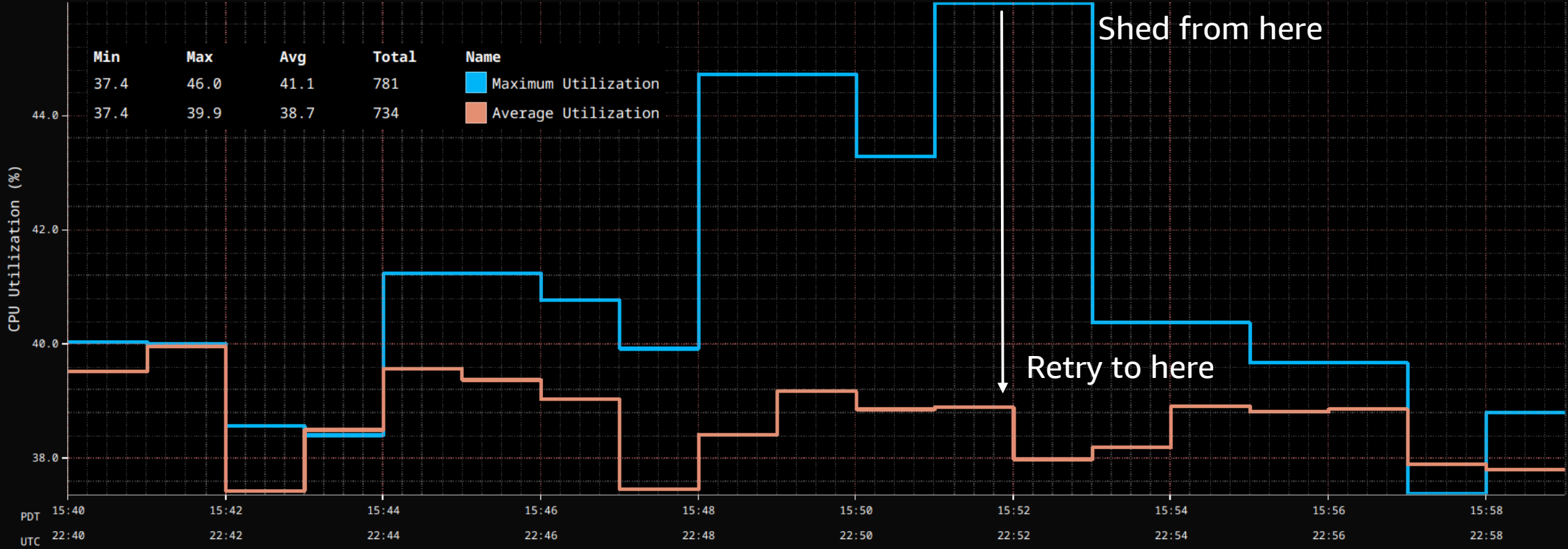
# Stay up – Prioritized load shedding



System Load with Buffer

4x

4x Failure Buffer

10%

2x

50%

2x Success Buffer

x

40%

System
Resources

Time

$$T_{transit} = 15m$$

# Shed the right load

# Are retries a good idea?



CPU Utilization Spread Max-Avg

| Min | Max | Avg | Total | Name |
|------|------|------|-------|------|
| 37.4 | 46.0 | 41.1 | 781 | Maximum Utilization |
| 37.4 | 39.9 | 38.7 | 734 | Average Utilization |

Shed from here

Retry to here

# Retry sparingly
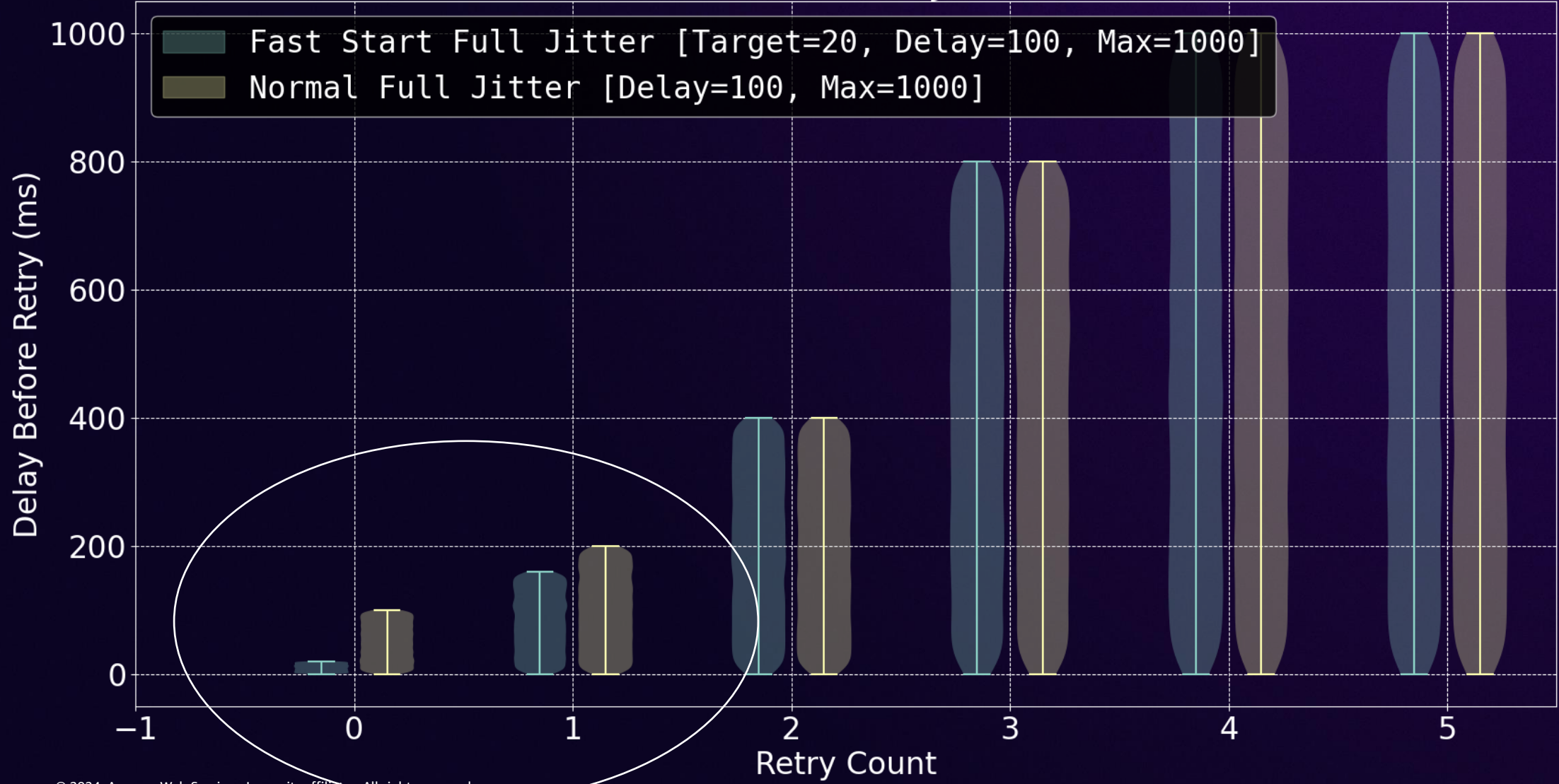
```
# Full jitter exponential backoff on shedding only
interceptor.retry.default.maxRetries     = 2
             i.r.d.statuses              = UNAVAILABLE
             i.r.d.backoffPolicy         = exponential
             i.r.d.backoffPolicy.jitterMode      = full
             i.r.d.backoffPolicy.targetMillis    = 20
             i.r.d.backoffPolicy.delayMillis     = 100
             i.r.d.backoffPolicy.maxDelayMillis = 1000
```

$$let\ R\ =\ \text{retry}\ \#\ \in\ [0,\ 1,\ 2,\ ...\ \text{retry}_{max}\ -\ 1]$$

$$base(R)\ =\ \min\left(\text{delay},\ \text{target}\ \times\ (R\ +\ 1)^2\right)$$

$$retry(R)\ =\ \text{rand}\left[0,\ \min\left\{\text{delay}_{max},\ base(R)\ \times\ 2^R\right\}\right)$$
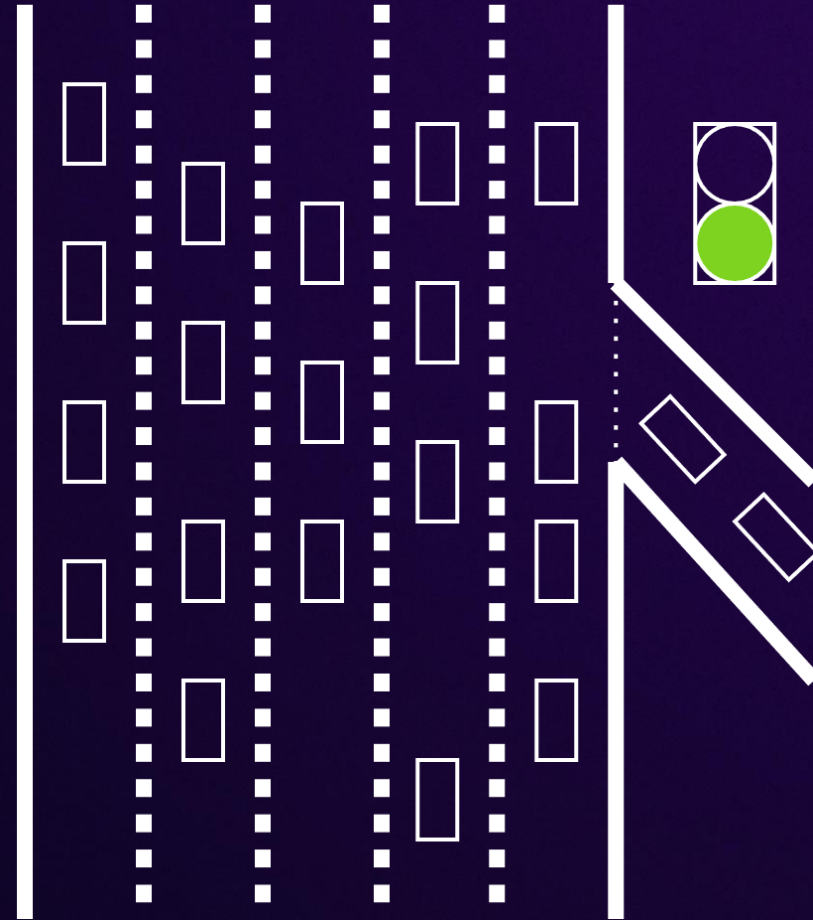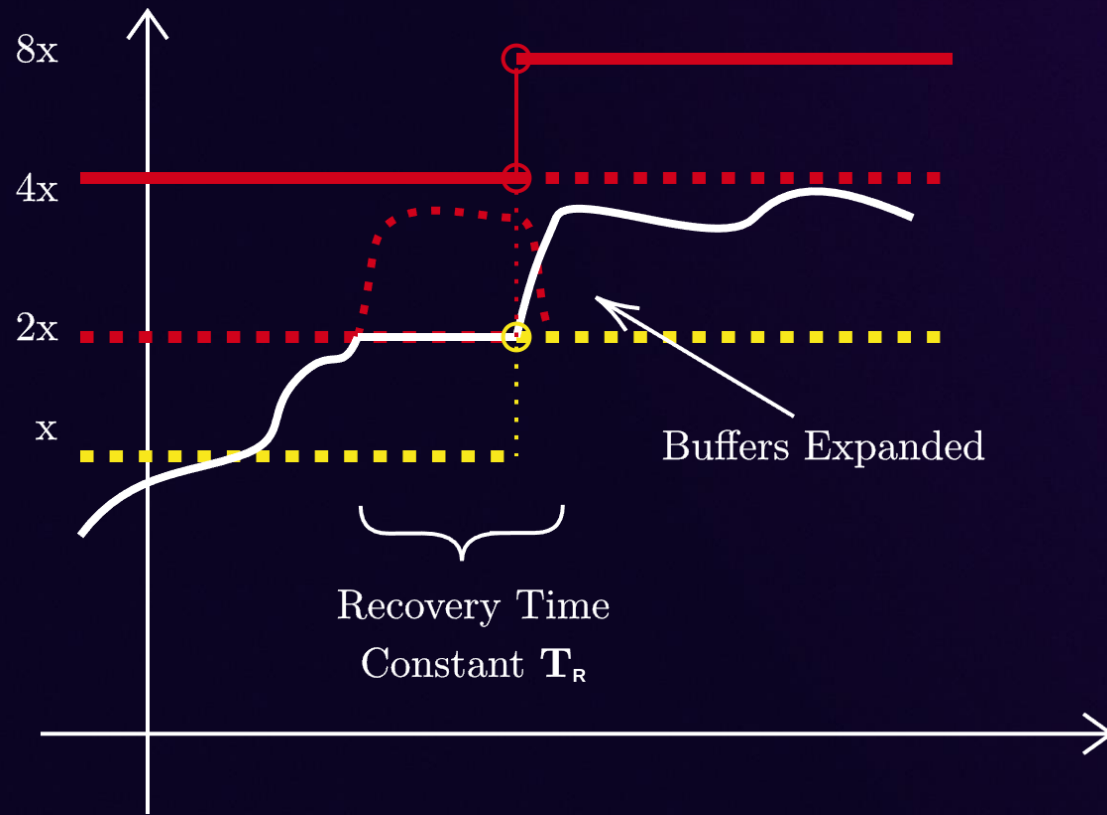
# Fast start full jitter



Fast Start Full Jitter

**Fast Start Full Jitter [Target=20, Delay=100, Max=1000]**
**Normal Full Jitter [Delay=100, Max=1000]**

Delay Before Retry (ms)

Retry Count

# More capacity is the real solution



Buffer Recovering After Load Spike

Buffers Expanded

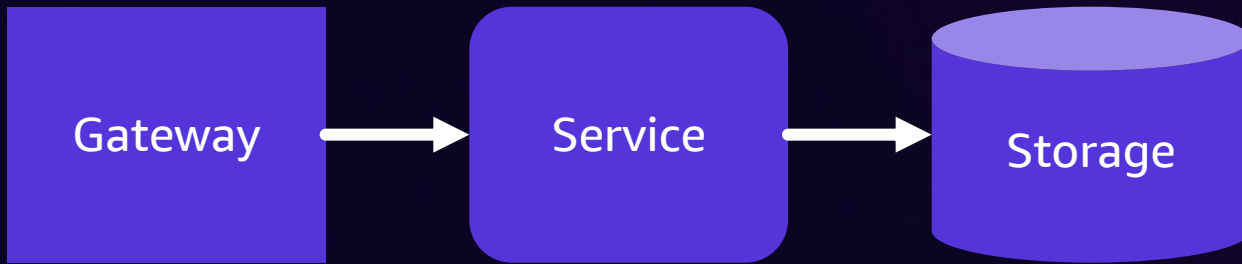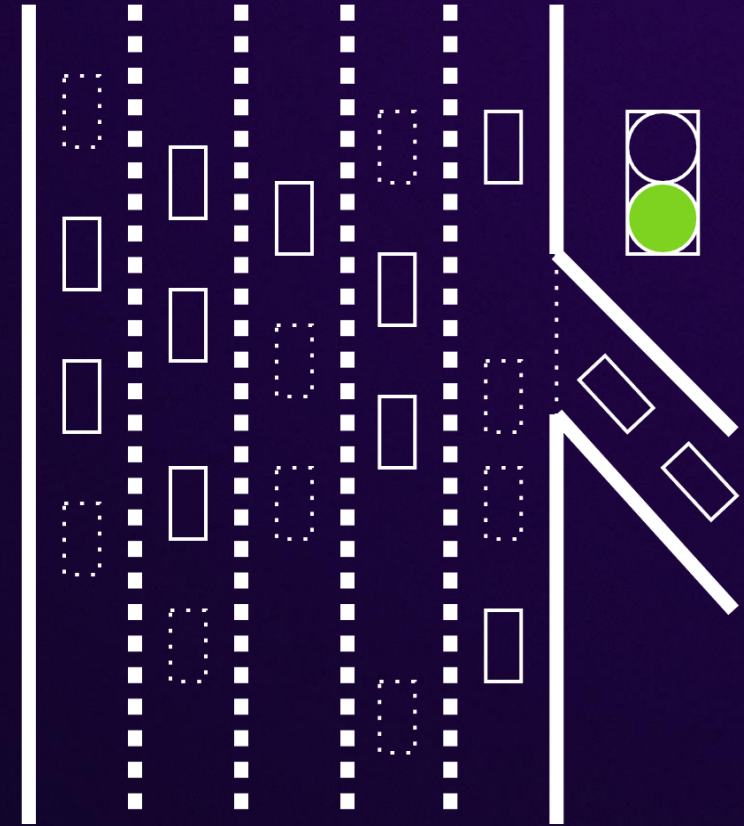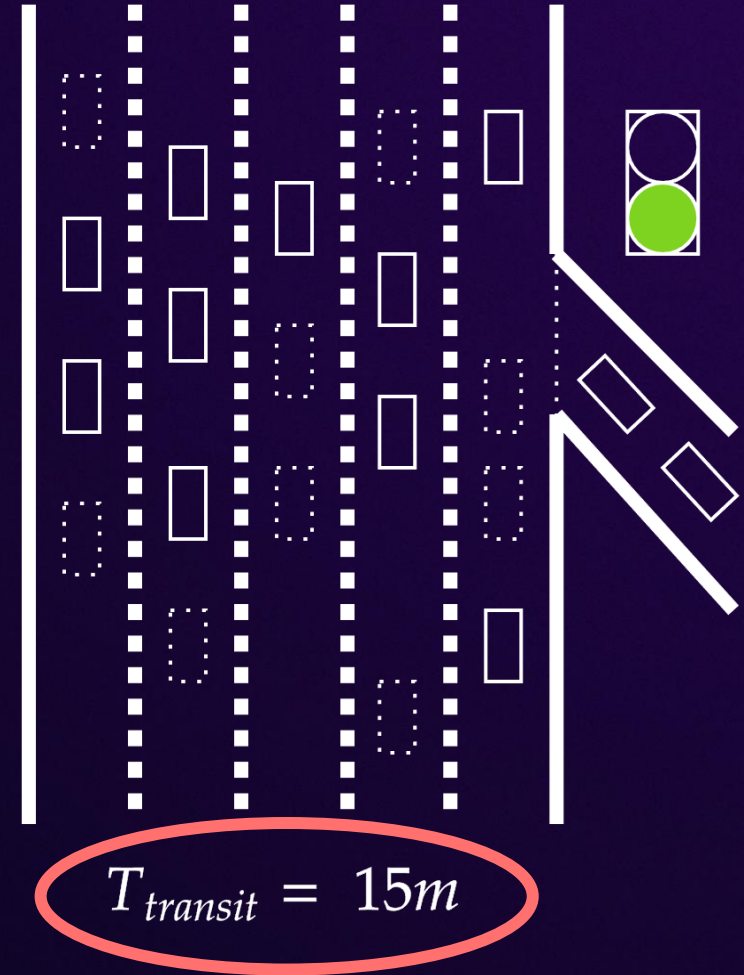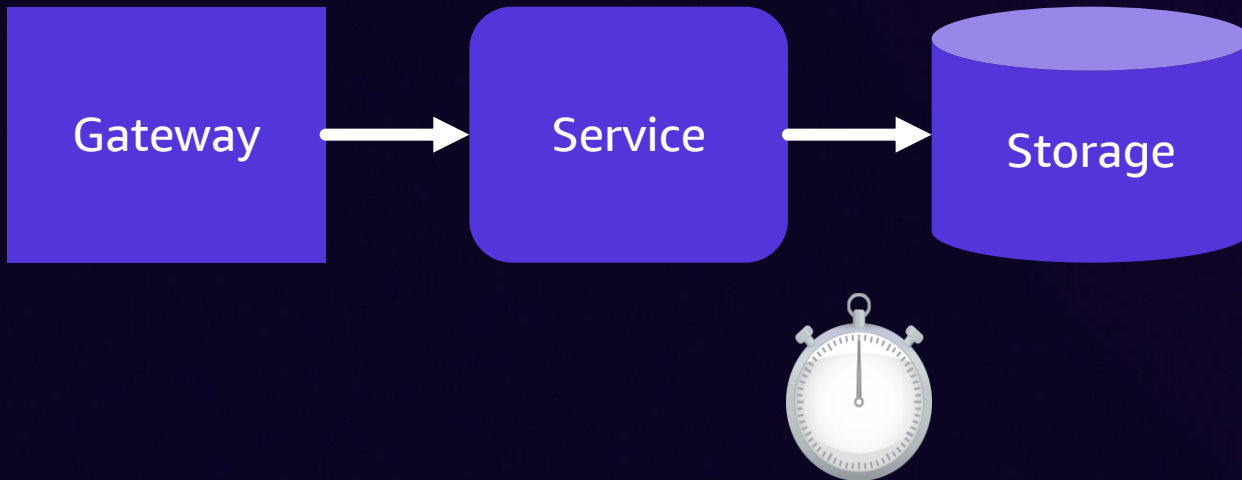Recovery Time Constant $T_R$

$$T_{transit} = 15m$$

# What about IO bound services?

CPU CAPACITY IS NECESSARY, BUT NOT SUFFICIENT!

Most services talk to other services

Async calls don't take much CPU

Gateway → Service → Storage

Latency???

$$T_{transit} = 15m$$

# What about IO bound services?

Most services talk to other services

Async calls don't take much CPU

Gateway → Service → Storage

$$T_{transit} = 15m$$

# Measure latency as utilization



Service Latency as Utilization

Legend:
- Normal Service Time μ=3, p95=240
- Target Latency @ 50ms

Target Latency Utilization = 0.21

X-axis: Request Latency (ms)
Y-axis: Requests Completed By (%)

# Measure latency as utilization



Service Latency as Utilization

Legend:
- Normal Service Time μ=3, p95=240
- Degraded Service Time μ=150, p95=1000
- Target Latency @ 50ms

Target Latency Utilization = 0.88!

Y-axis: Requests Completed By (%)
X-axis: Request Latency (ms)

# Stay up – Allocate IO success buffer



```
resource:
    utilization:
        kv-slo:
            target: 40
            max: 80
    limiter:
        kv-slo:
            enabled: true
            utilization:
                source: kv-slo
            buffer: success
```

# Stay up – Add IO limiters

# Stay up – Add prioritized IO limiters
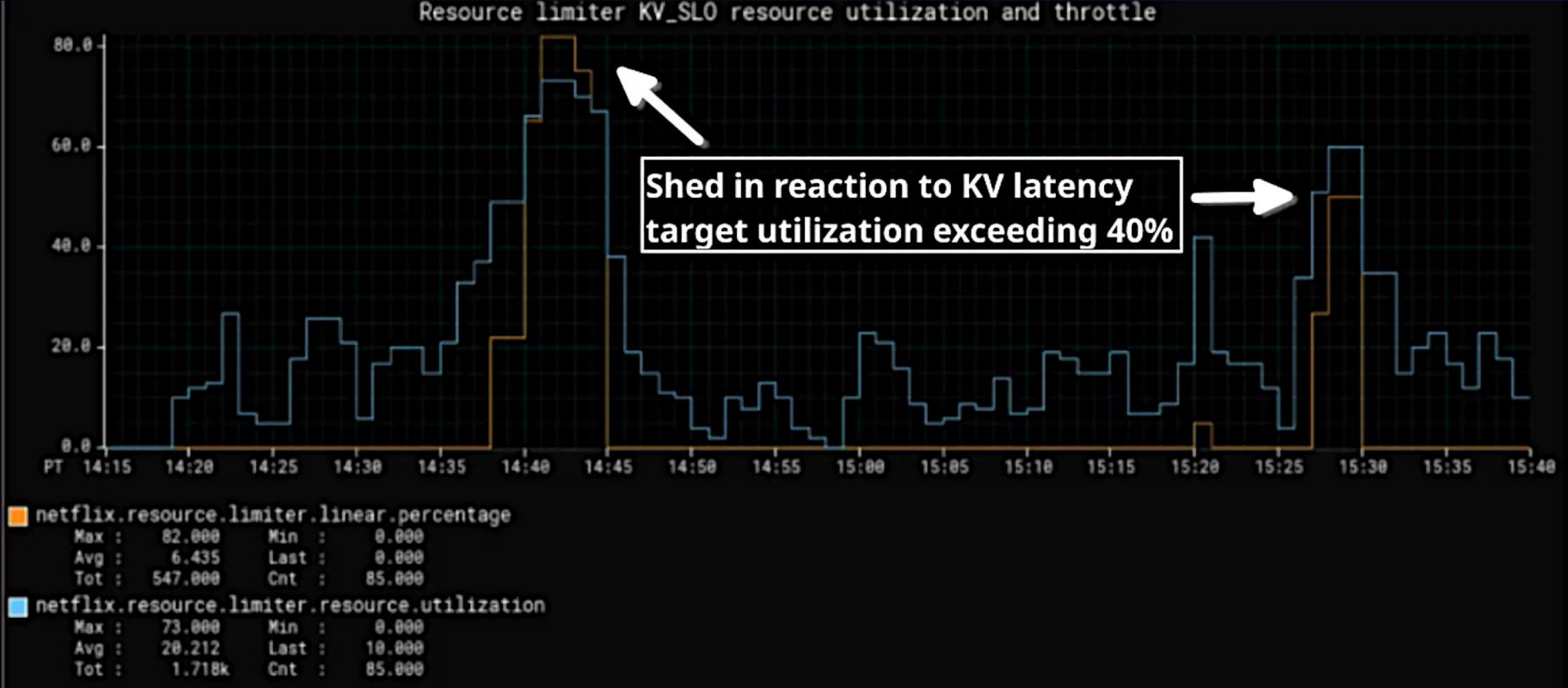


Never shed
- High-priority writes

At 80% *max* shed
- High-priority reads

At 40% *target* utilization shed
- Low-priority reads

# Stay up – Add IO limiters



Resource limiter KV_SLO resource utilization and throttle

Shed in reaction to KV latency target utilization exceeding 40%

netflix.resource.limiter.linear.percentage
| | | | |
|---|---|---|---|
| Max : | 82.000 | Min : | 0.000 |
| Avg : | 6.435 | Last : | 0.000 |
| Tot : | 547.000 | Cnt : | 85.000 |

netflix.resource.limiter.resource.utilization
| | | | |
|---|---|---|---|
| Max : | 73.000 | Min : | 0.000 |
| Avg : | 20.212 | Last : | 10.000 |
| Tot : | 1.718k | Cnt : | 85.000 |

# Stay up – Prioritized shedding

## Success buffer shedding

*Prioritized* **[CPU]**

*Prioritized* **[Latency target]**

## Failure buffer shedding

*Unprioritized* **[CPU]**

**[Latency timeout]**
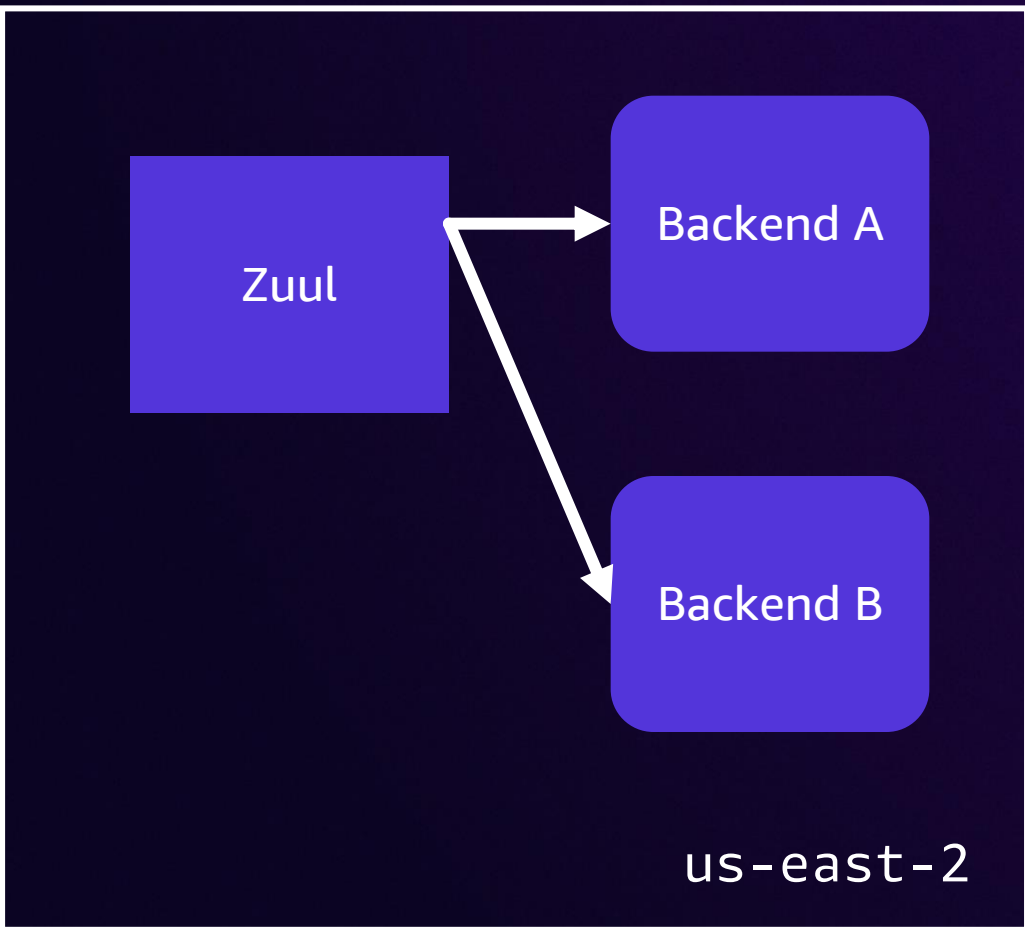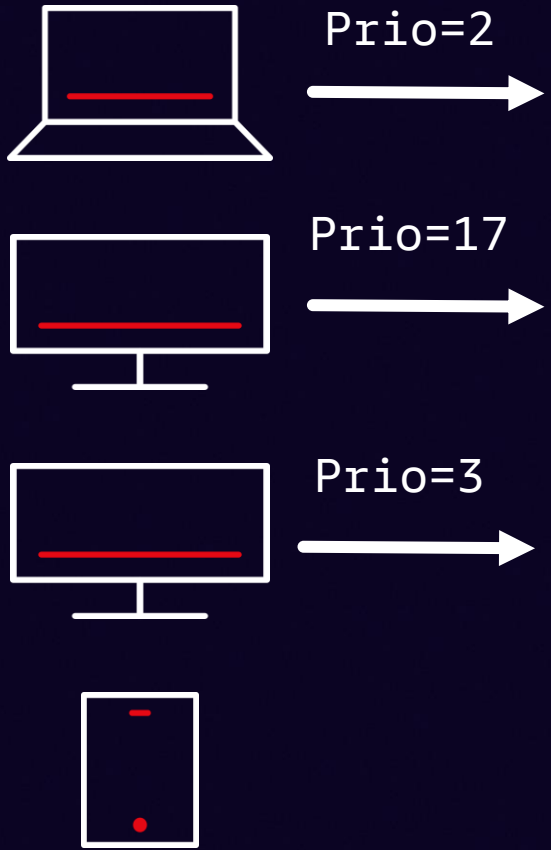
### Generic IO based load-shedding

Some services are not CPU-bound but instead are IO-bound by backing services or datastores that can apply back pressure via increased latency when they are overloaded either in compute or in storage capacity. For these services we re-use the prioritized load shedding techniques, but we introduce new utilization measures to feed into the shedding logic. Our initial implementation supports two forms of latency based shedding in addition to standard adaptive concurrency limiters (themselves a measure of average latency):

1. The service can specify per-endpoint target and maximum latencies, which allow the service to shed when the service is abnormally slow regardless of backend.

2. The Netflix storage services running on the Data Gateway return observed storage target and max latency SLO utilization, allowing services to shed when they overload their allocated storage capacity.

These utilization measures provide early warning signs that a service is generating too much load to a backend, and allow it to shed low priority work before it overwhelms that backend. The main advantage of these techniques over concurrency limits alone is they require less tuning as our services already must maintain tight latency service-level-objectives (SLOs), for example a p50 < 10ms and p100 < 500ms. So, rephrasing these existing SLOs as utilizations allows us to shed low priority work early to prevent further latency impact to high priority work. At the same time, the system *will accept as much work as it can* while maintaining SLO's.

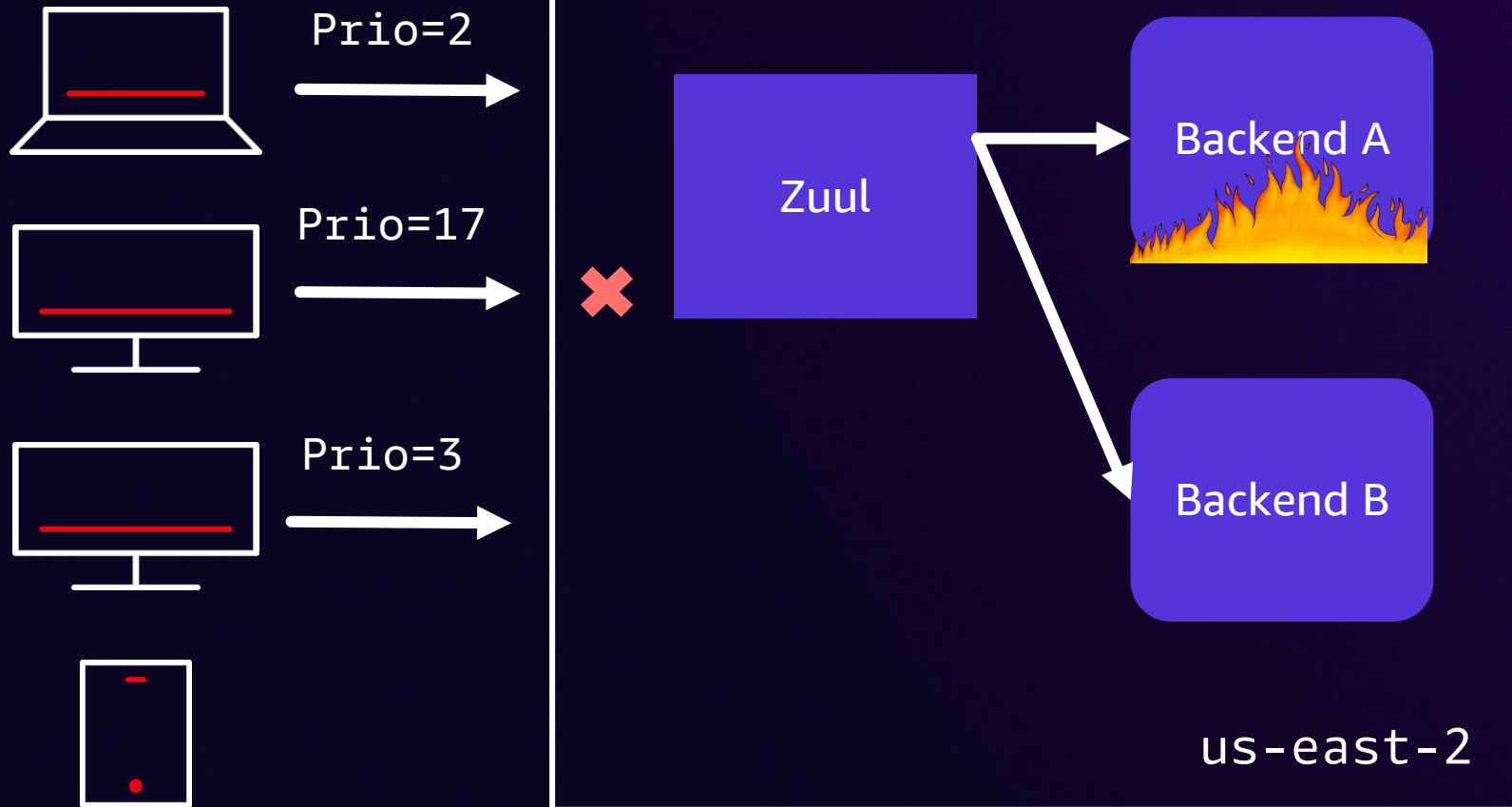https://netflixtechblog.com/enhancing-netflix-reliability-with-service-level-prioritized-load-shedding-e735e6ce8f7d

# Fallbacks!



Prio=2

Prio=17

Prio=3

Zuul → Backend A, Backend B

us-east-2

https://netflixtechblog.com/keeping-netflix-reliable-using-prioritized-load-shedding-6cc827b02f94

# Fallbacks!



Prio=2

Prio=17

Prio=3

Zuul

Backend A

Backend B

us-east-2

# Fallback and shift



Prio=2

Prio=17

Zuul

Backend A

Prio=99

Zuul

Prio=3

cassandra

cassandra

us-east-2

Full
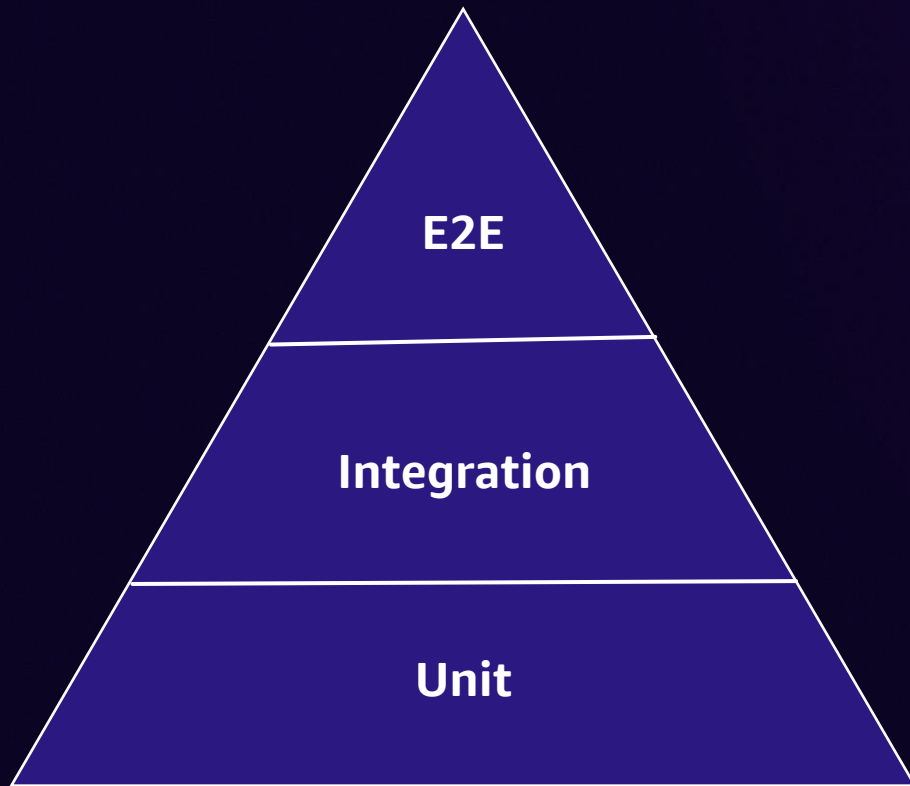Active

us-east-1

# Fallback and shift
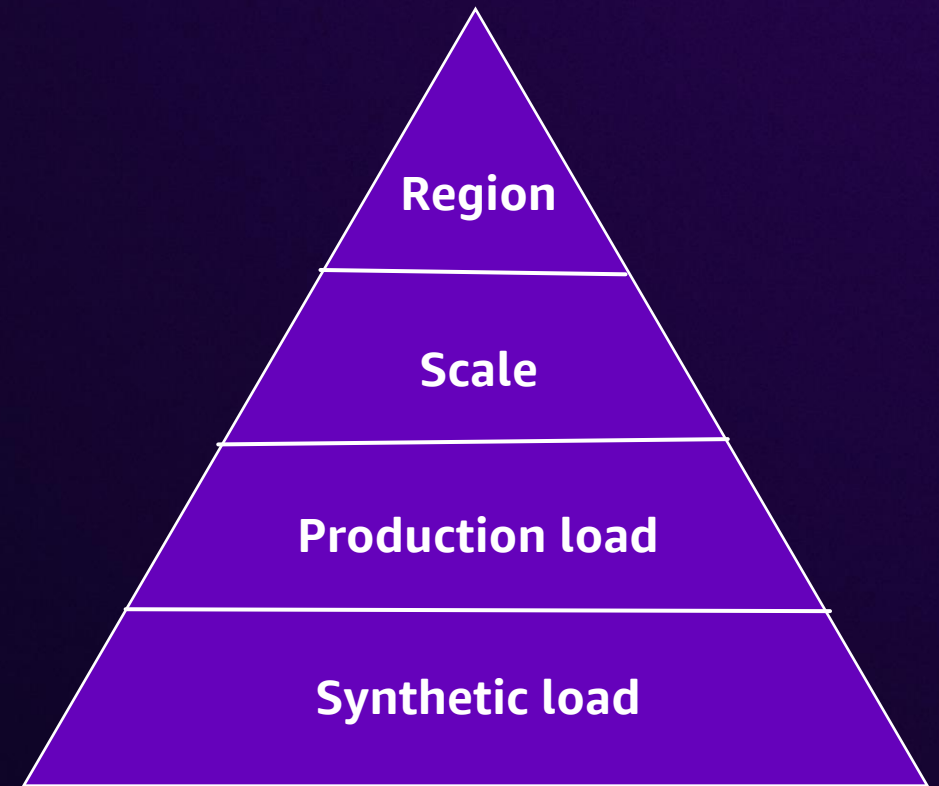
# 05: Resilience testing
## Validating the techniques

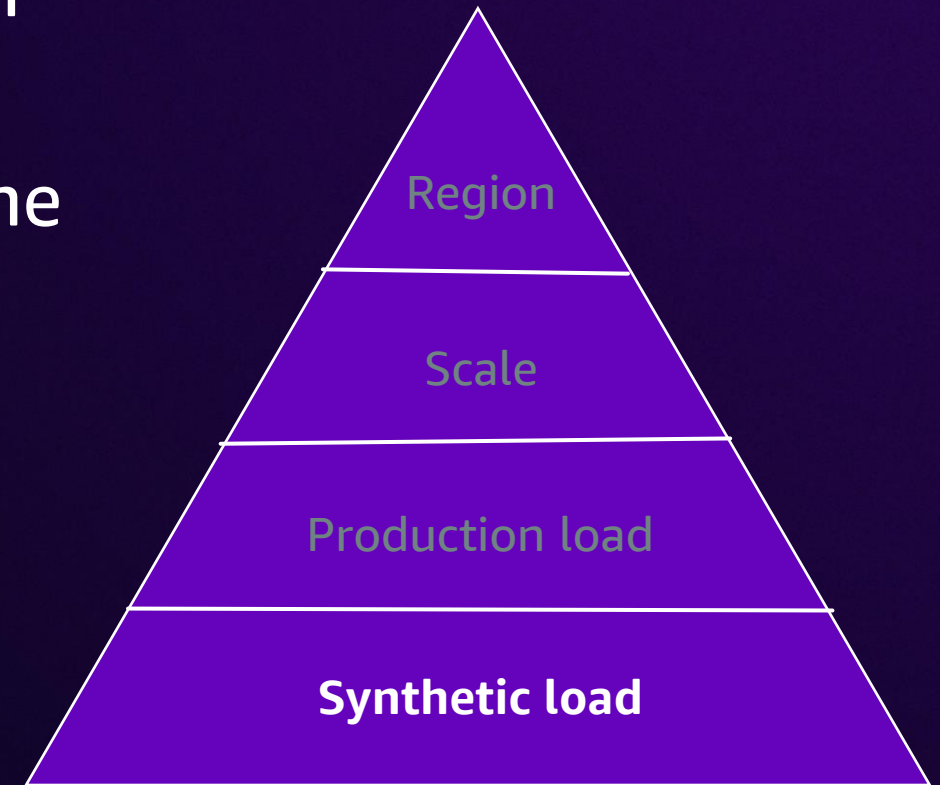# The resilience testing pyramid
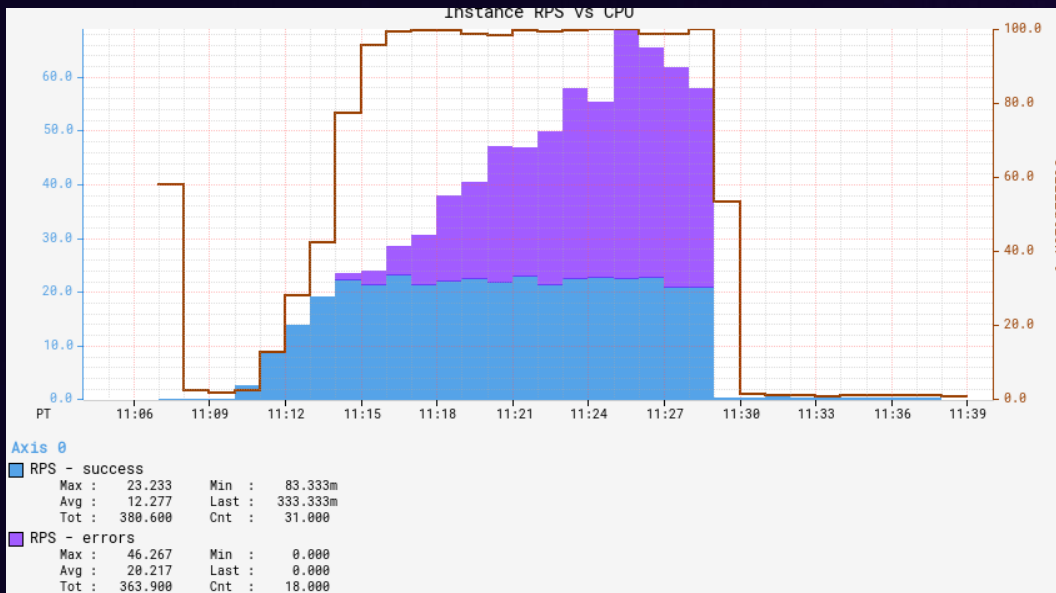
Testing pyramid

Resilience testing pyramid

# Service-level synthetic load testing

Use synthetic traffic to test an application in isolation

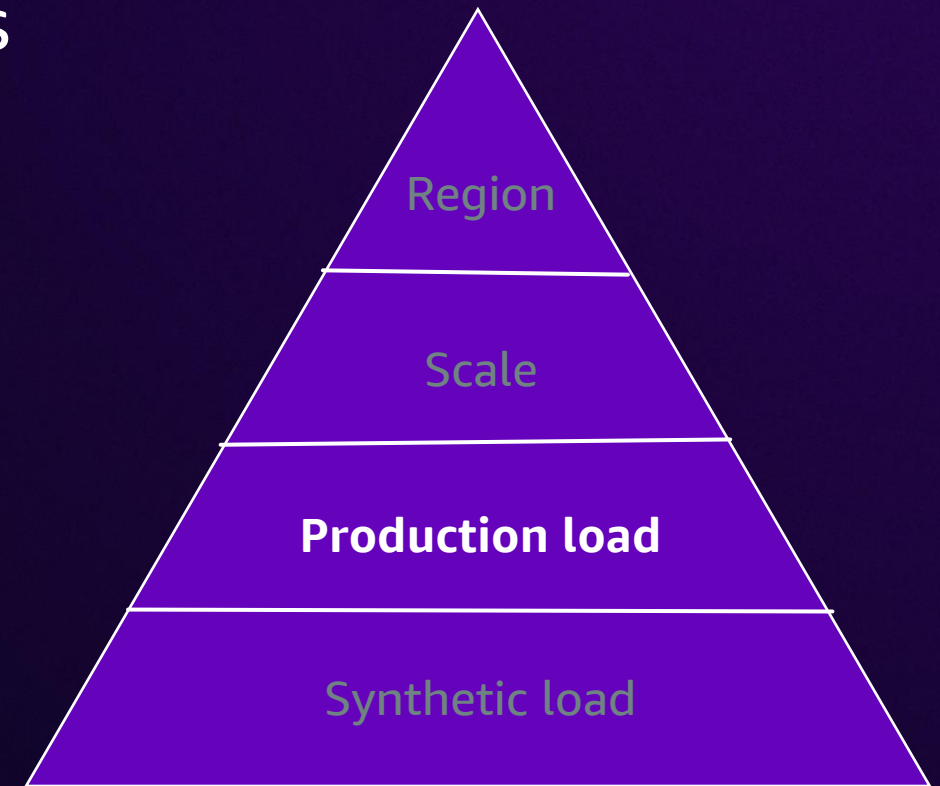Find bottlenecks in application code and tune load shedding configs

# Production load testing

Autoscaling squeeze test through our Chaos Automation Platform

Introduces a load spike to a service to test how load shedding and autoscaling behave

Tests the actual production config with real traffic

Region

Scale
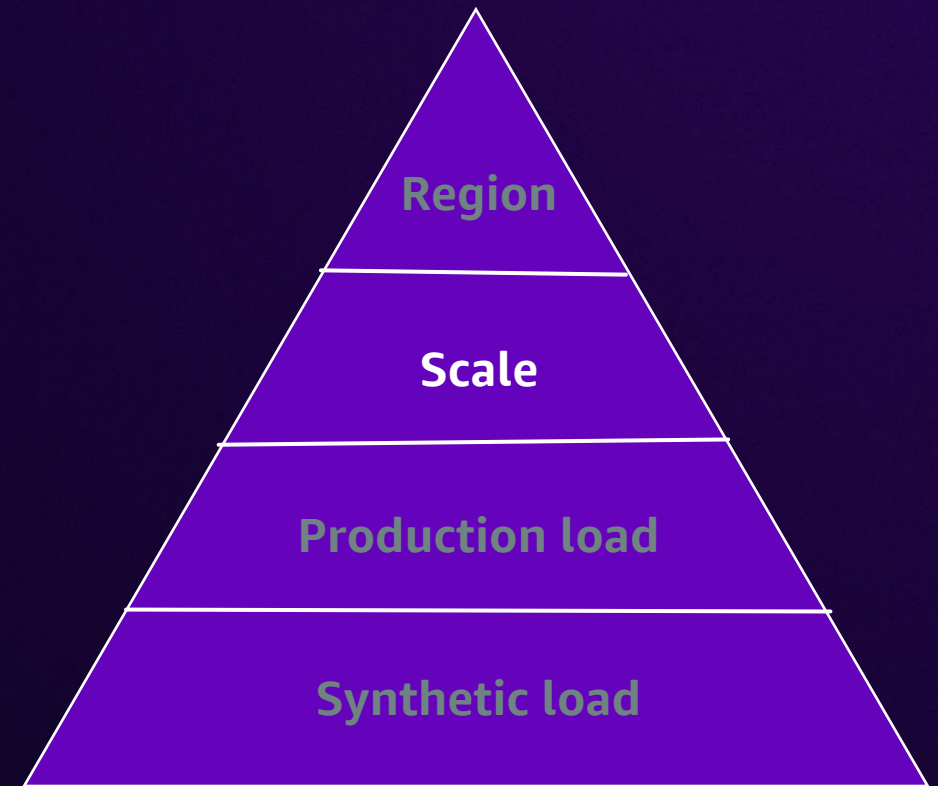
**Production load**

Synthetic load

# Region scale

Move all global traffic into a single region

Uses regional failover tooling

Finds issues only seen at scale: load that scales with:
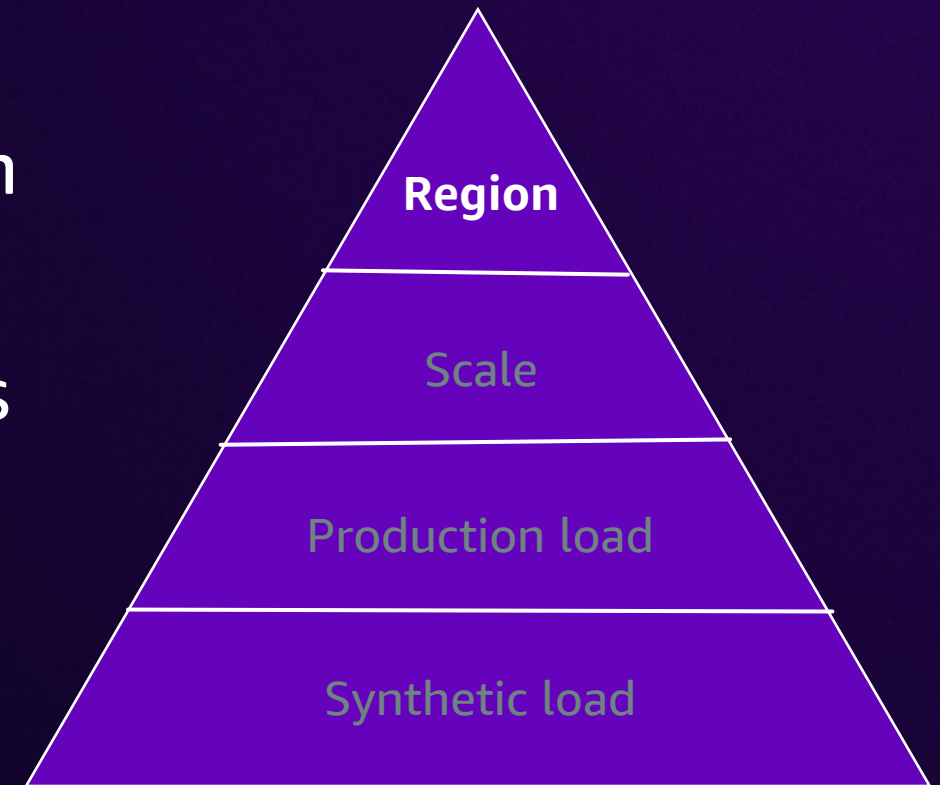
- # of instances

- # of RPS



aws

# Region load testing

E2E tests that simulate user behavior

Uses synthetic traffic against the production
Netflix API

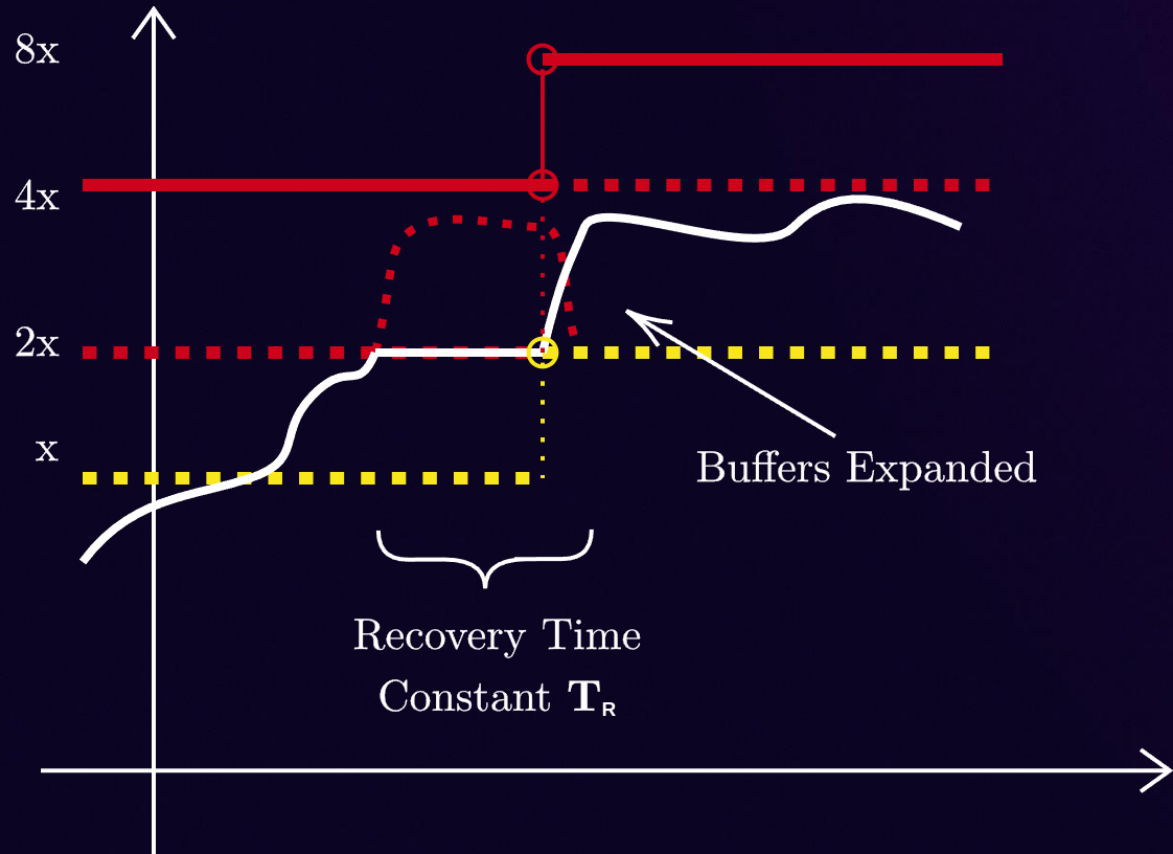Simulate title launches and failure scenarios

Lets us test scales that are even bigger than
our current global peak

**Region**

Scale

Production load

Synthetic load

# 06: Conclusions and wrap-up

# Measuring success



Buffer Recovering After Load Spike

**Goals:**

- Reduce time to recover

- Use regional failover less as the primary remediation

- Build resiliency assuming load spikes are the norm

# Takeaways

Handling load spikes is a mix of proactive and reactive mechanisms: investing in both is important!

Use your existing compute resources to answer only the most important requests. Fail quickly when overloaded.

Test. In Prod. As often as possible.

# Thank you!

**Please complete the session survey in the mobile app**

**Rob Gulewich**

rgulewich@netflix.com

**Ryan Schroeder**

rschroeder@netflix.com

**Joseph Lynch**

josephl@netflix.com

jolynch.github.io/