

# Techniques Netflix Uses to Weather Significant Demand Shifts

Joseph Lynch

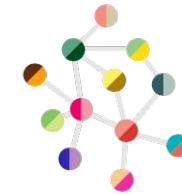
## Speaker



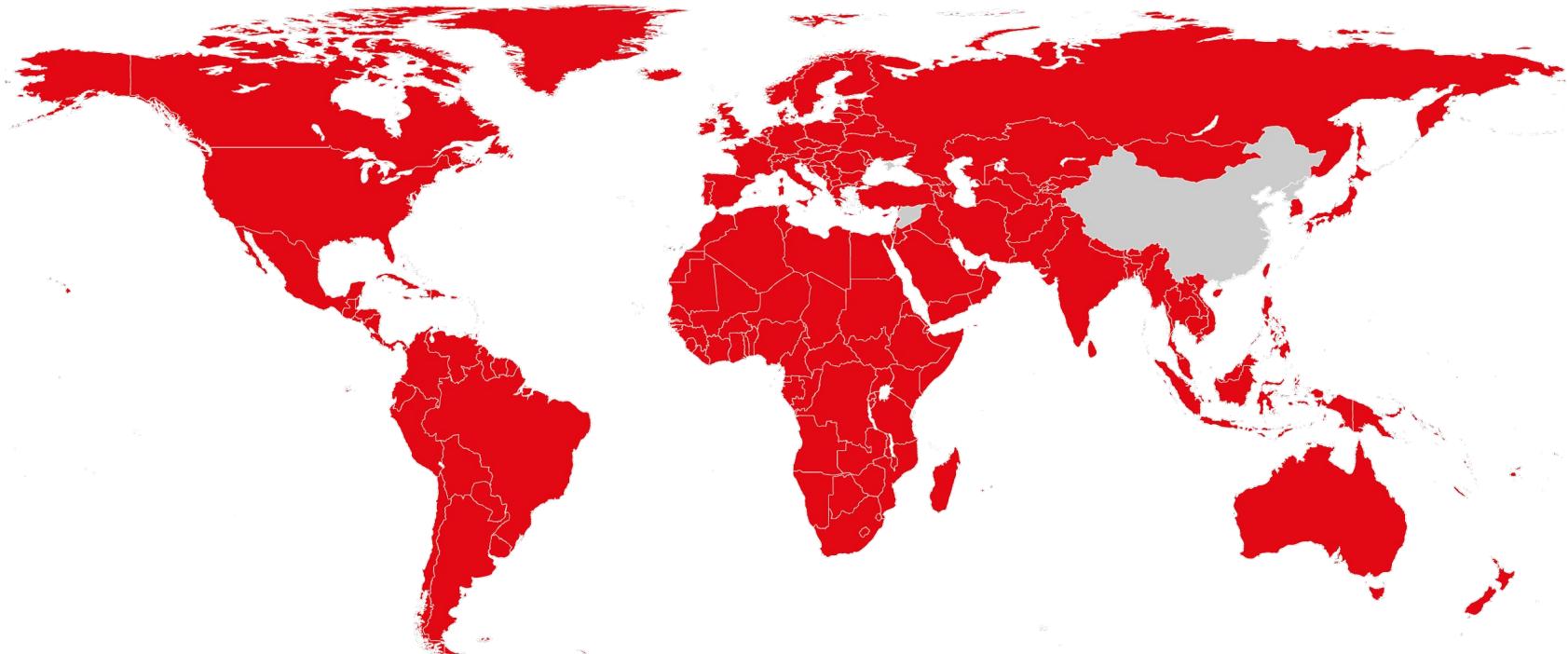
# Joseph Lynch

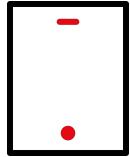
Principal Software Engineer  
Platform Engineering at Netflix

Database shepherd, compute  
optimizer, distributed system nerd



# Problem - Global



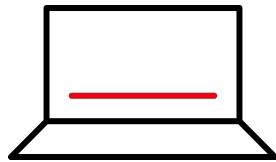


## Mobile

Diverse Device Capabilities

Weaker Network

Android and IOs

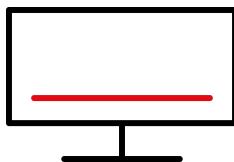


## PC

Diverse Device Capabilities

More Stable Network

Medium screens



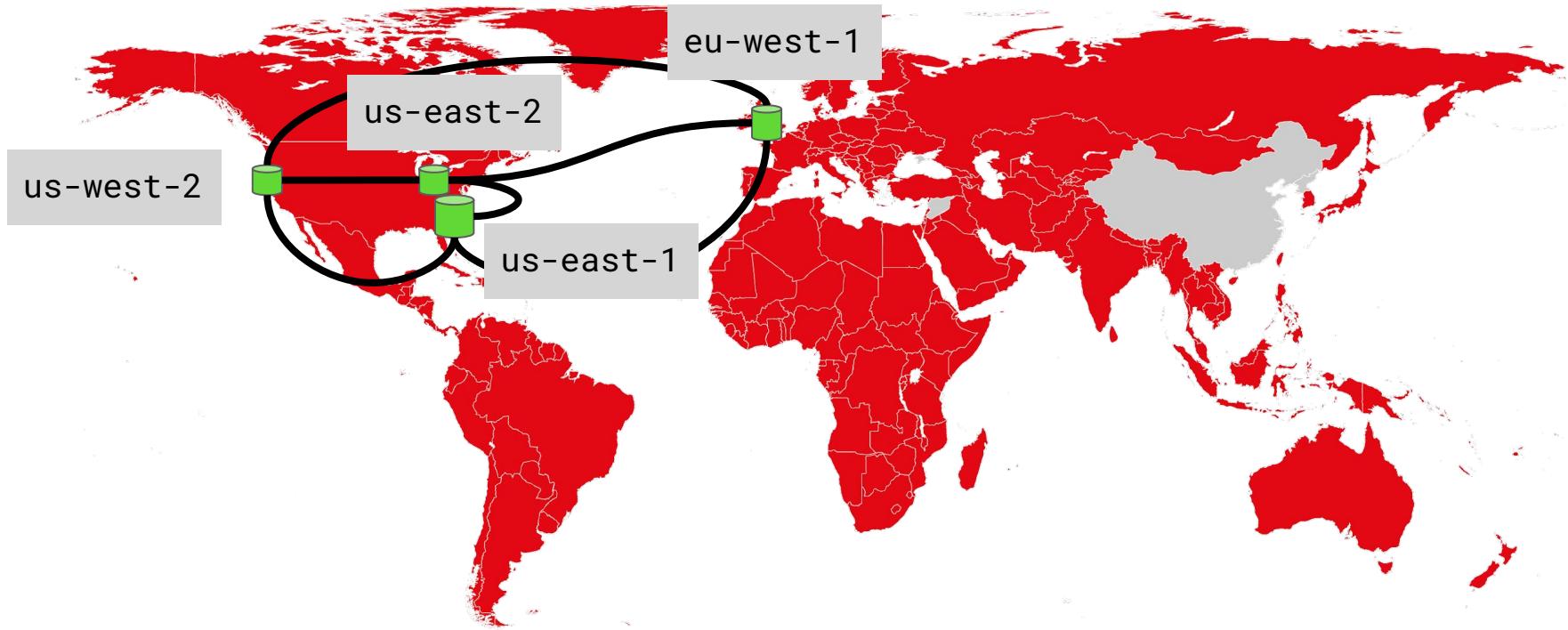
## TV

Diverse Device Capabilities

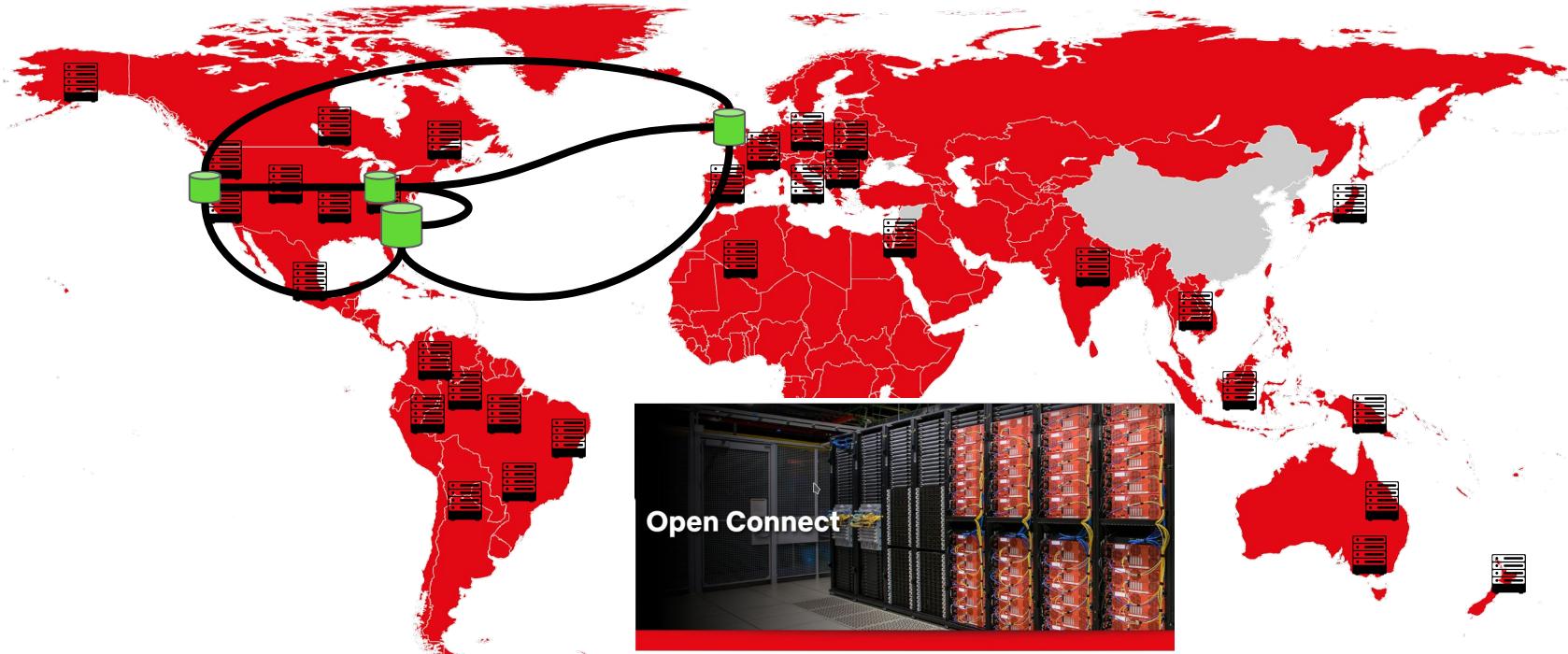
More Stable Network

Large screen

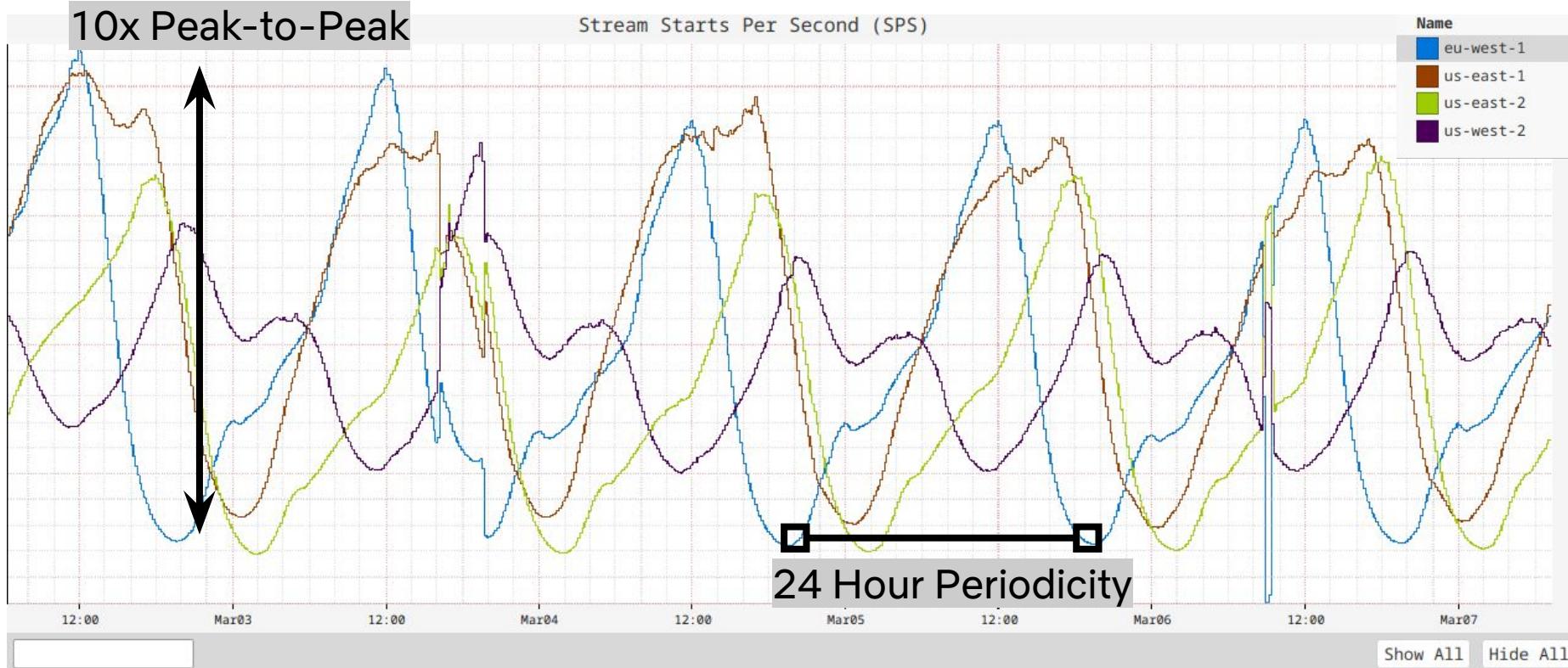
# Solution - Global Control Plane



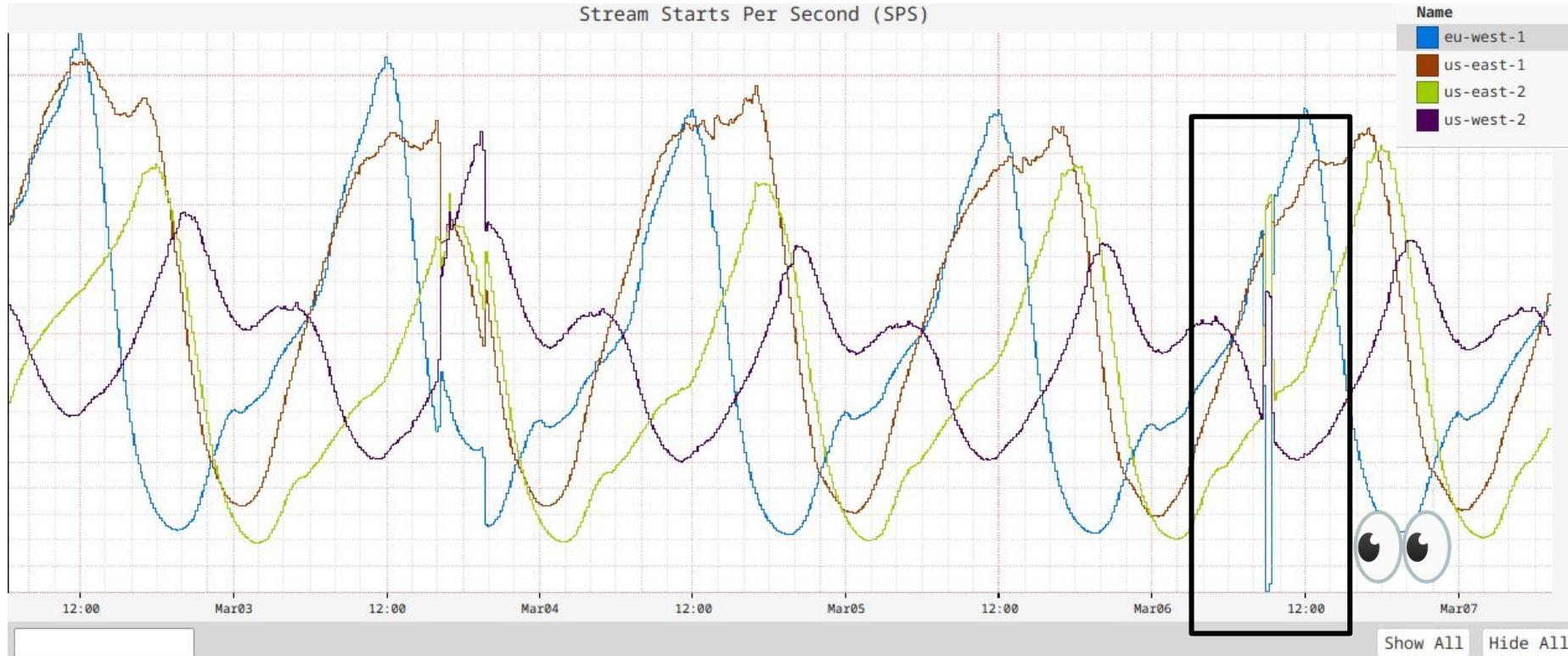
# Solution - Global Data Plane



# Variable Start-Per-Second (SPS) Load

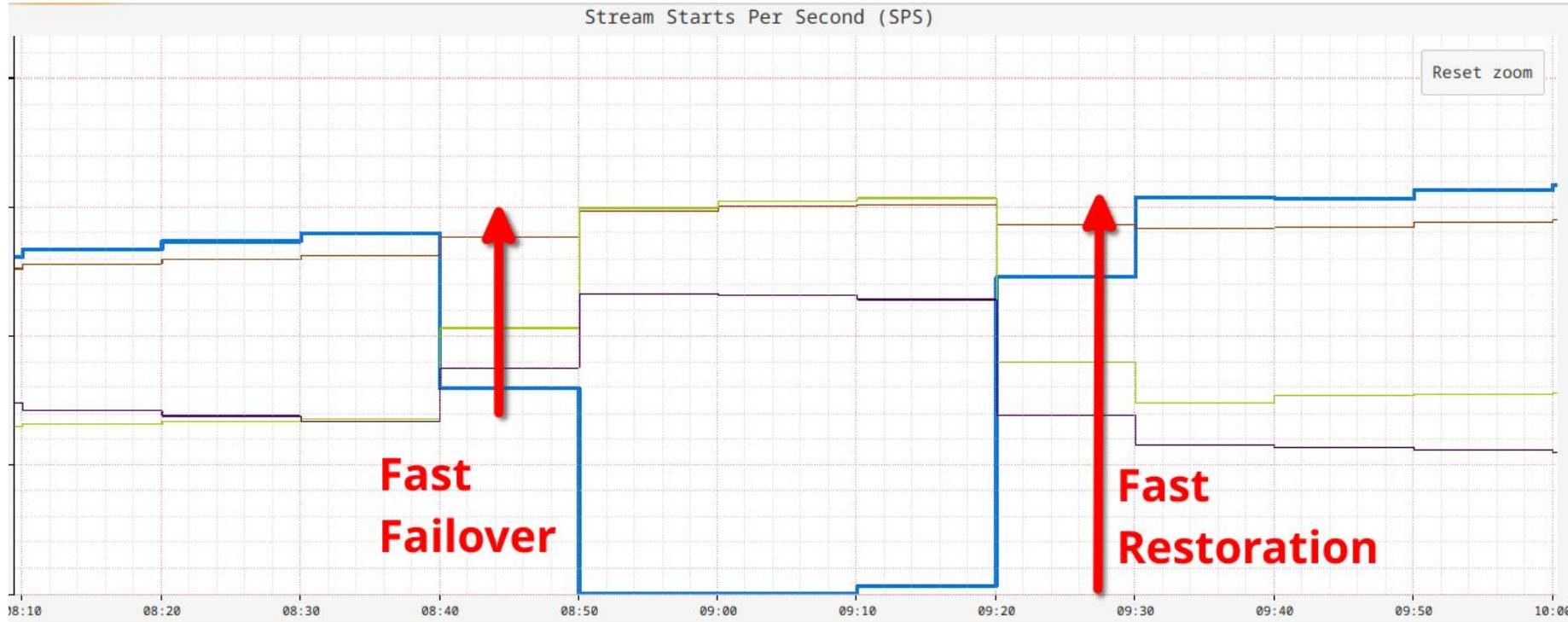


# Variable Start-Per-Second (SPS) Load

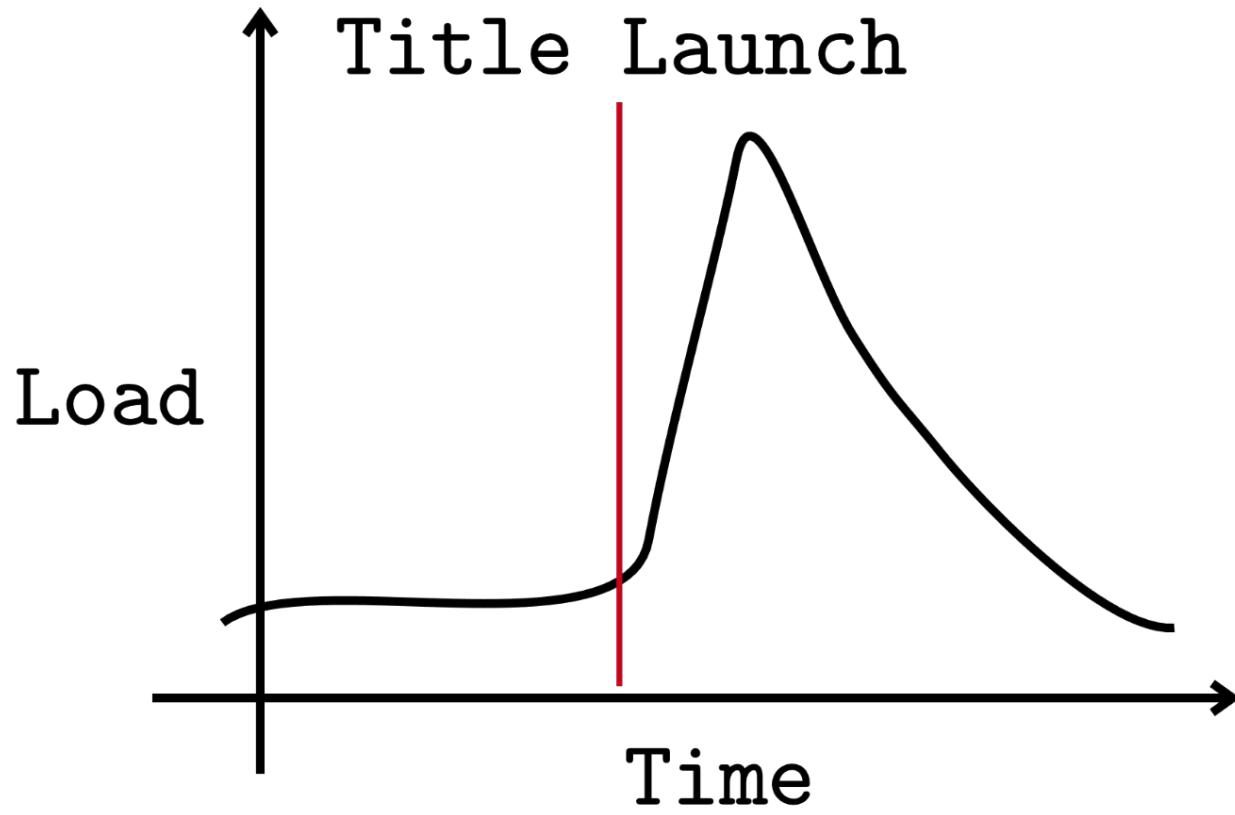


N

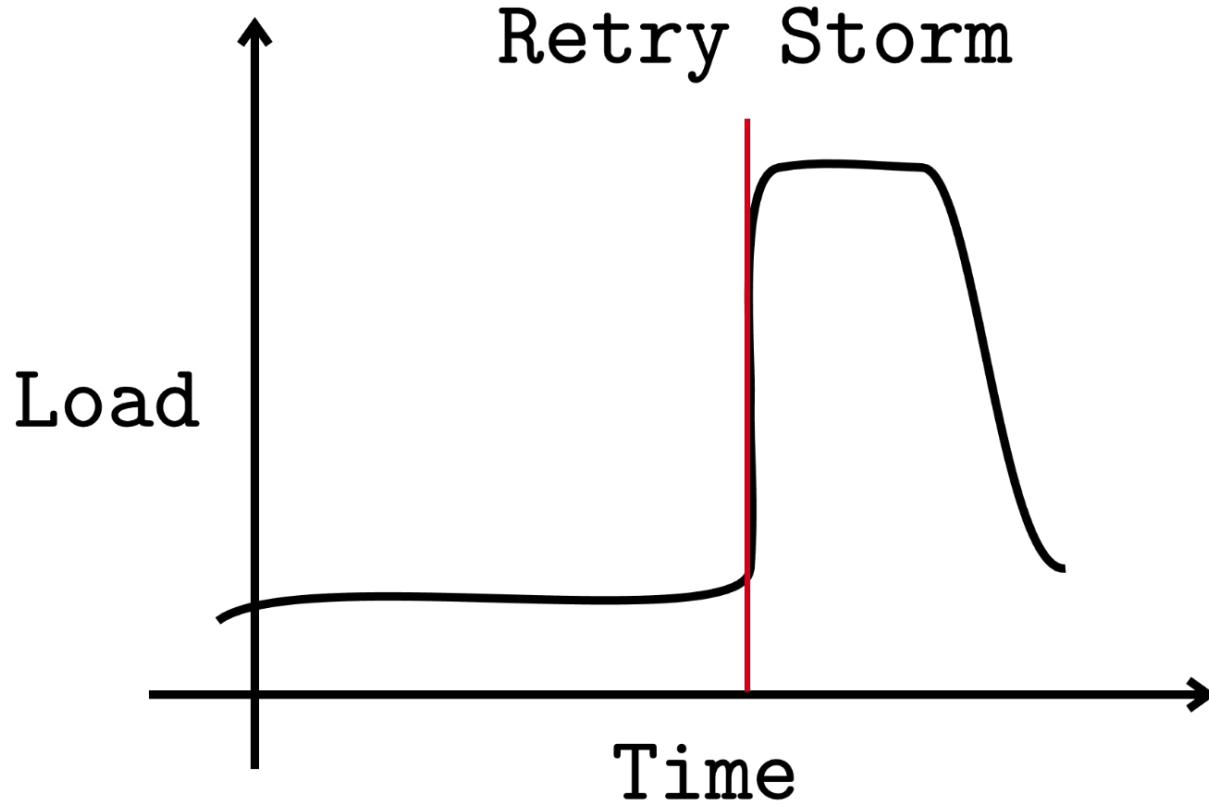
# Failover Driven Demand



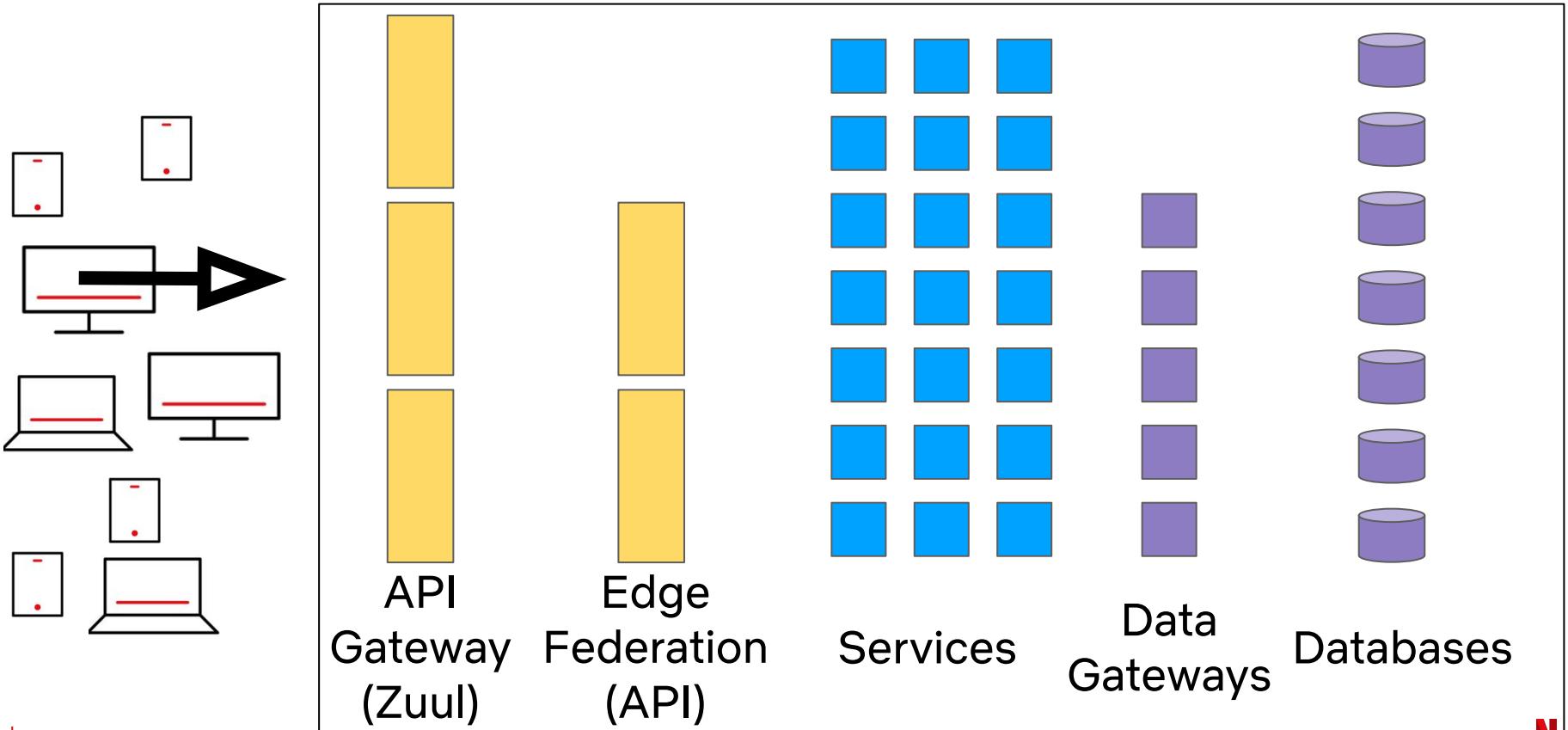
## Content Driven Demand



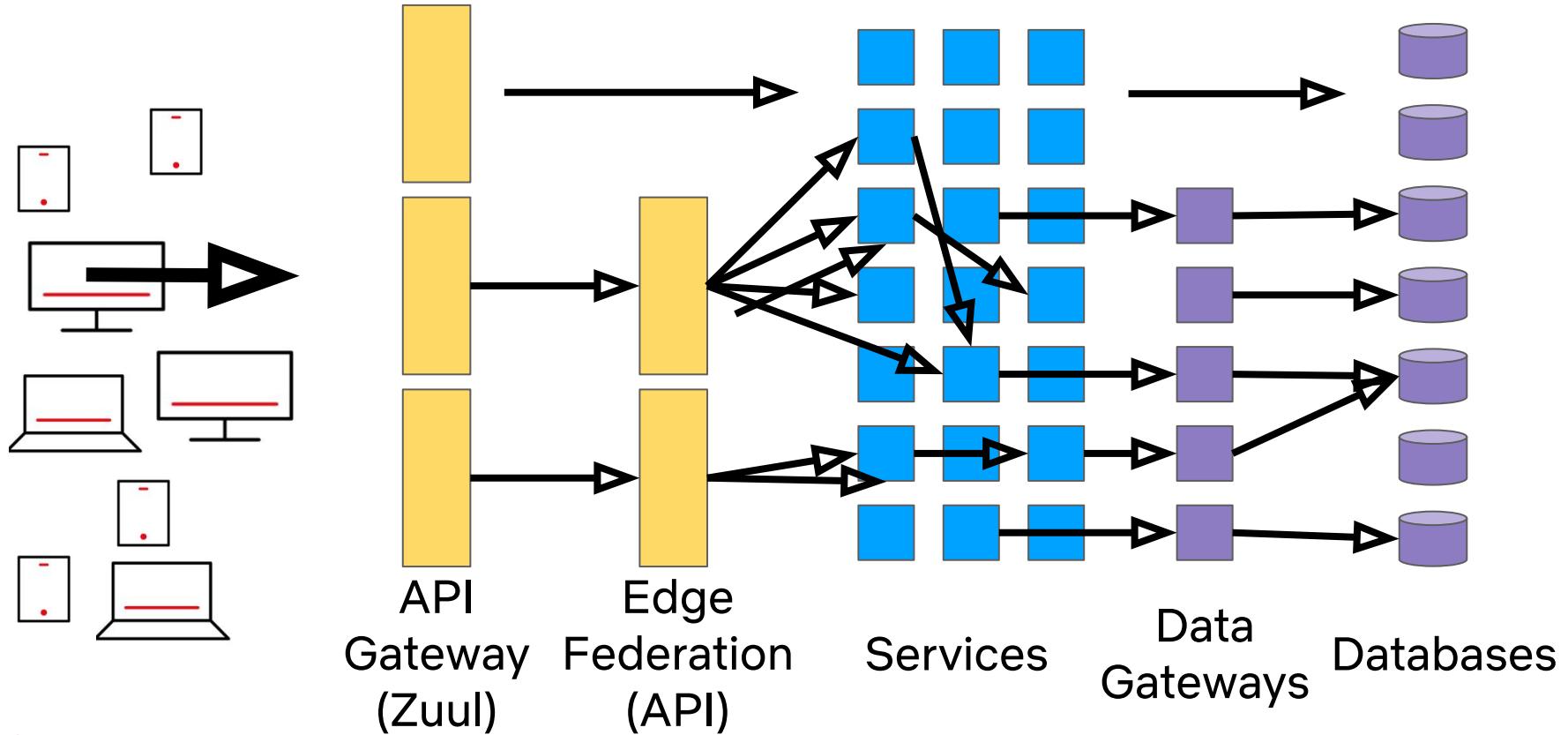
## Device Driven Demand



# Microservices



# Microservices





N

Traffic **Demand**

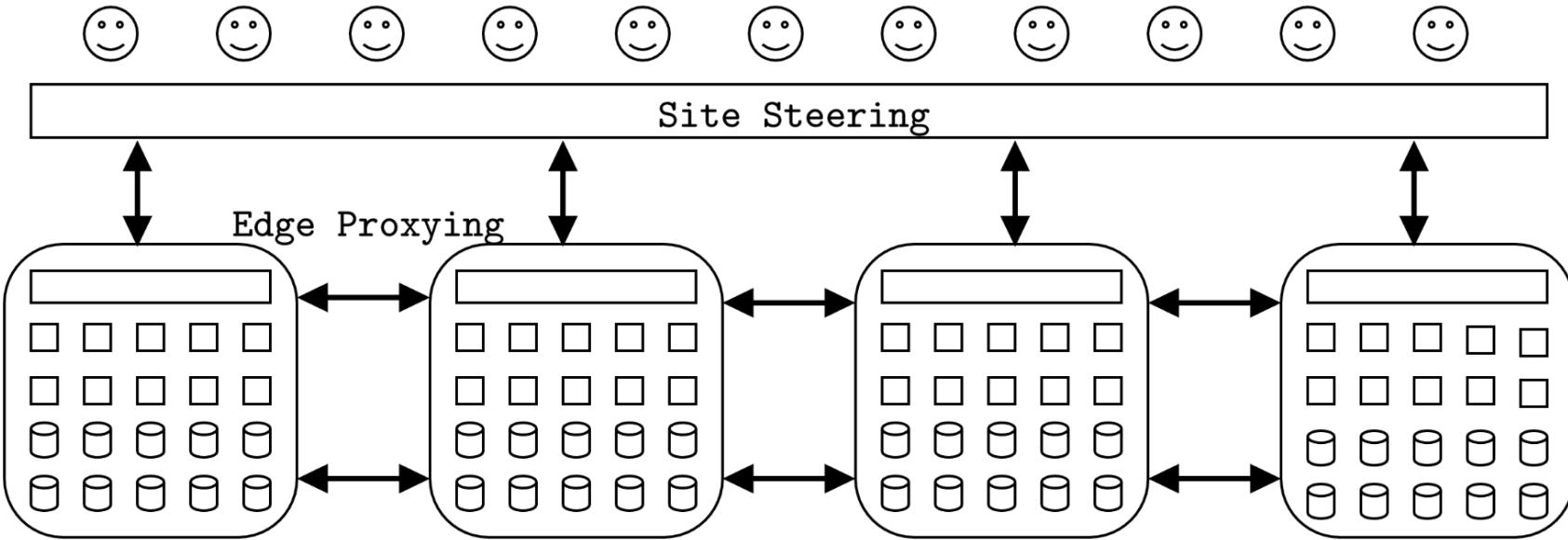
Compute **Supply**

**Resilience**  
Techniques

# Global Traffic Demand

Global architecture  
Traffic balancing  
Traffic shifting

# Global Architecture



Full Active Database and Cache Replication



# Latency

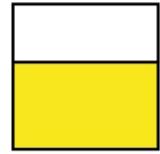
Bias users towards the **closest**  
Region or PoP

# Availability

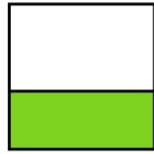
Bias users towards **least loaded**  
Region or PoP

# Control Plane Steering

us-west-2



us-east-2



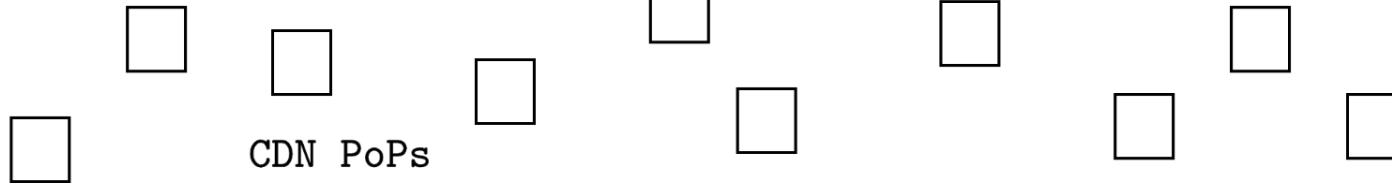
us-east-1



eu-west-1



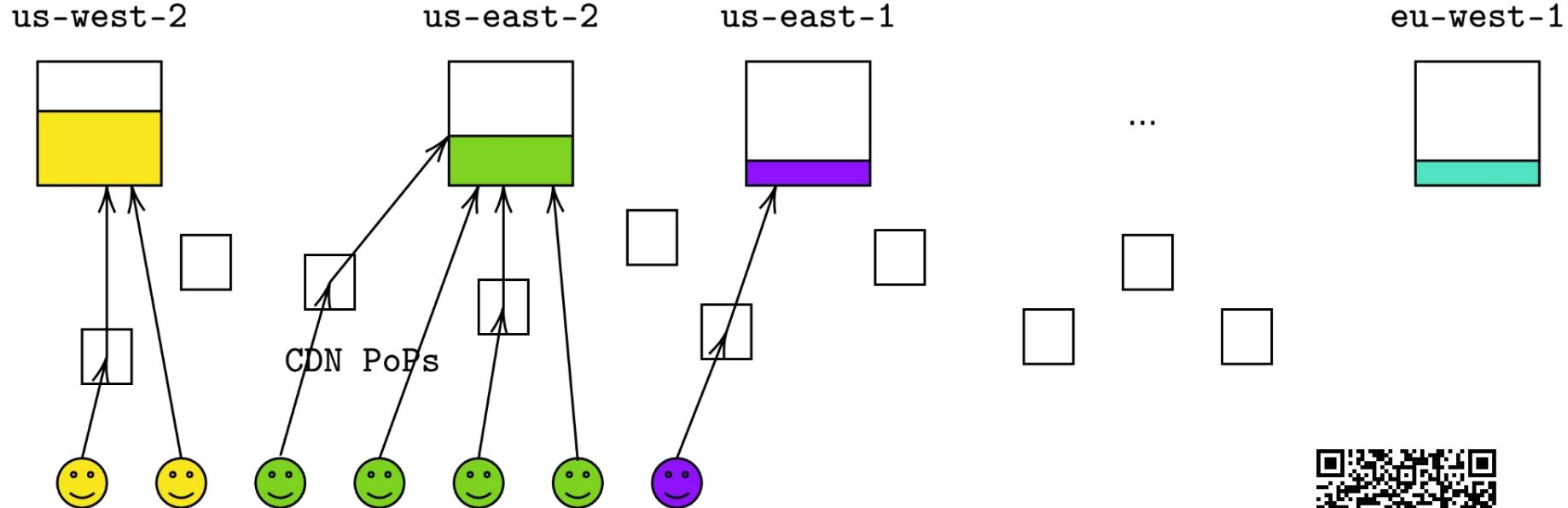
...



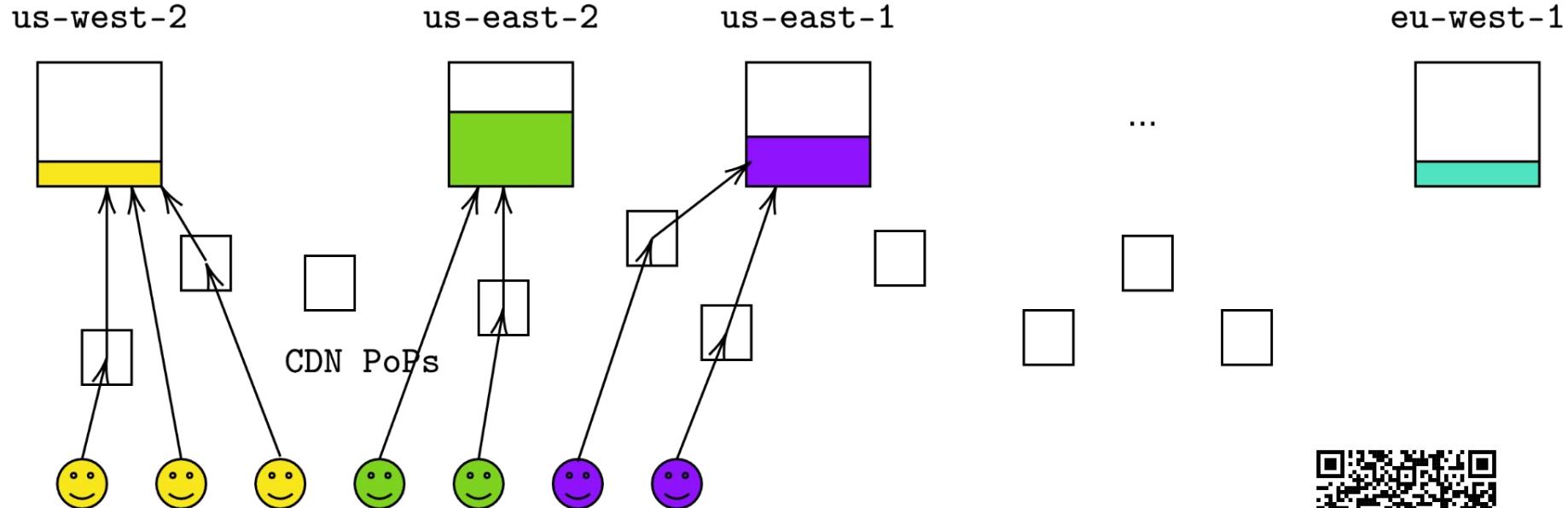
CDN PoPs



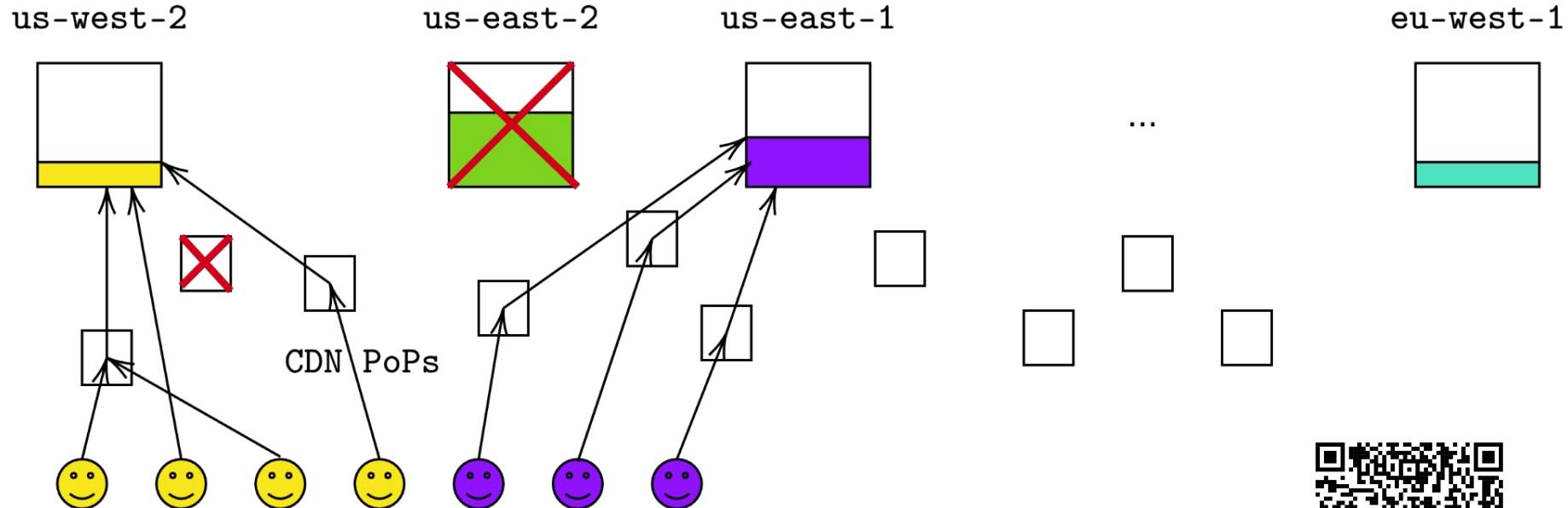
# Control Plane Steering



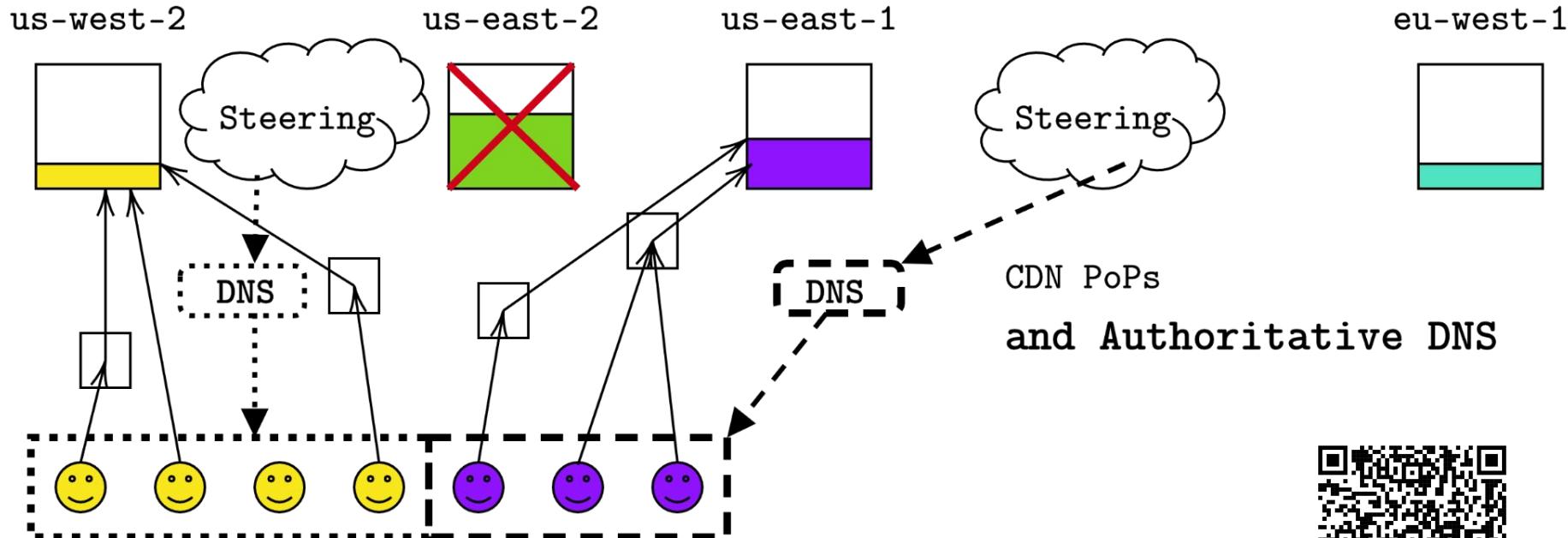
# Control Plane Steering



# Control Plane Steering



# Solution: Netflix DNS and Steering



# Predict

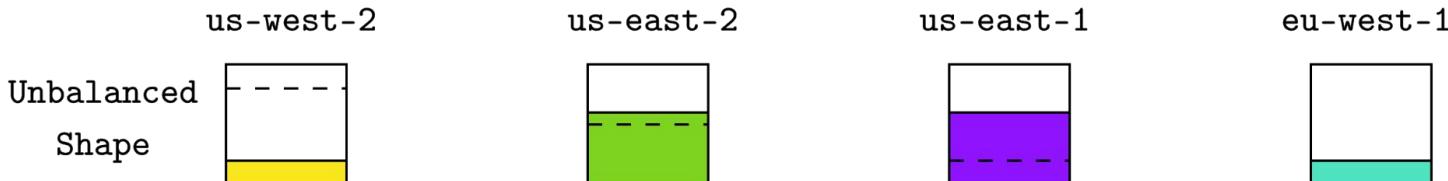
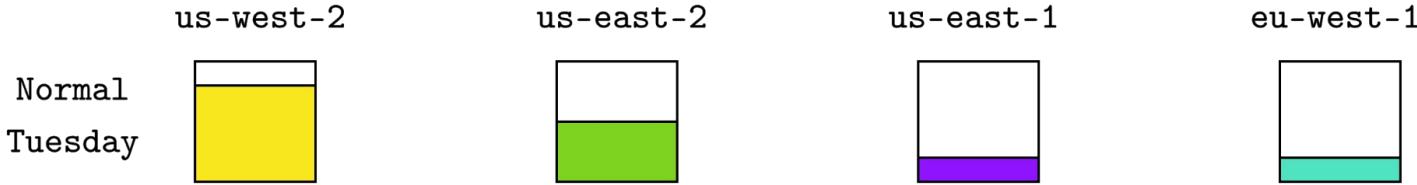
Shape traffic in anticipation

# React

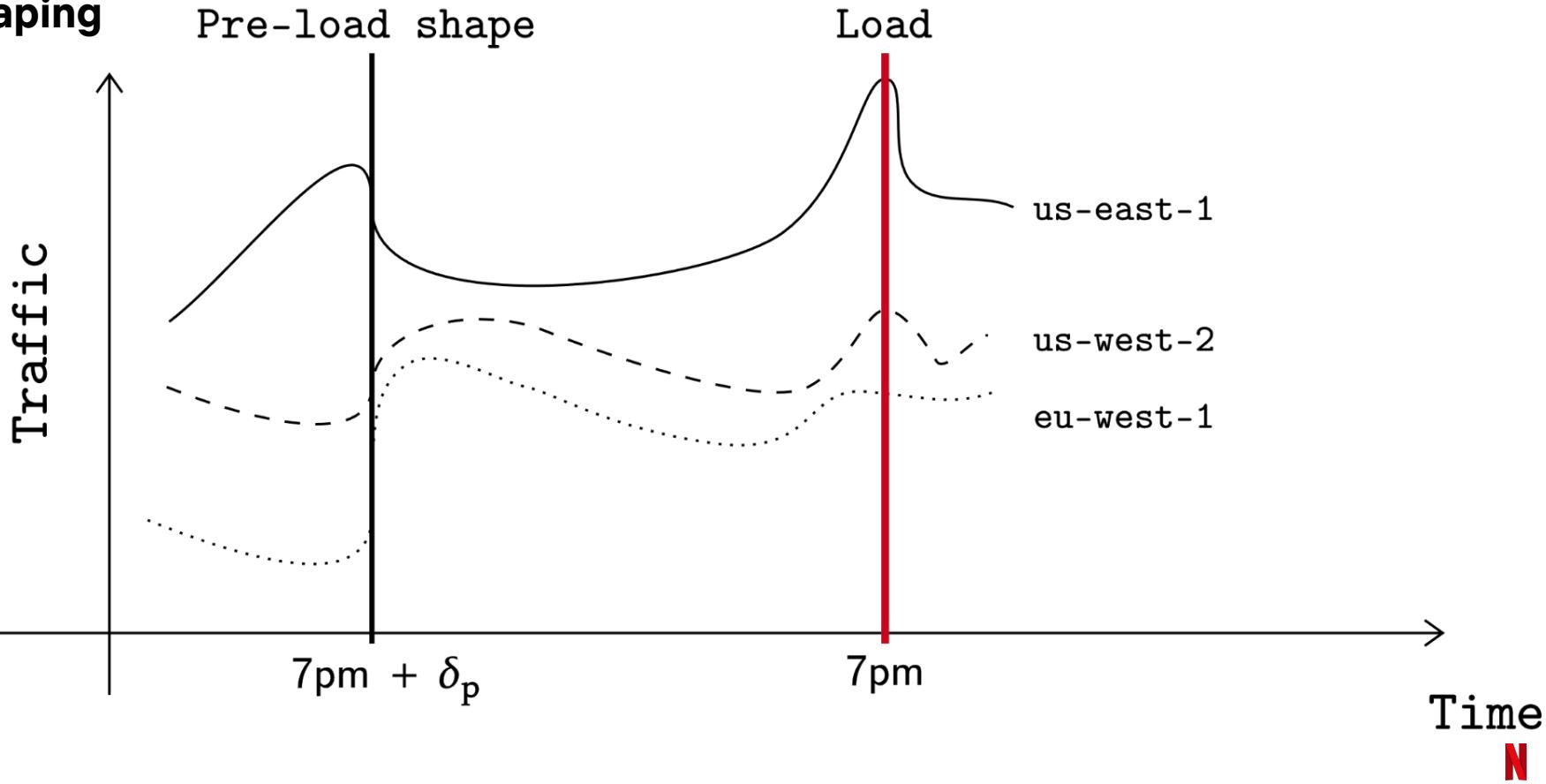
Restore balance

## Predict Traffic Spikes

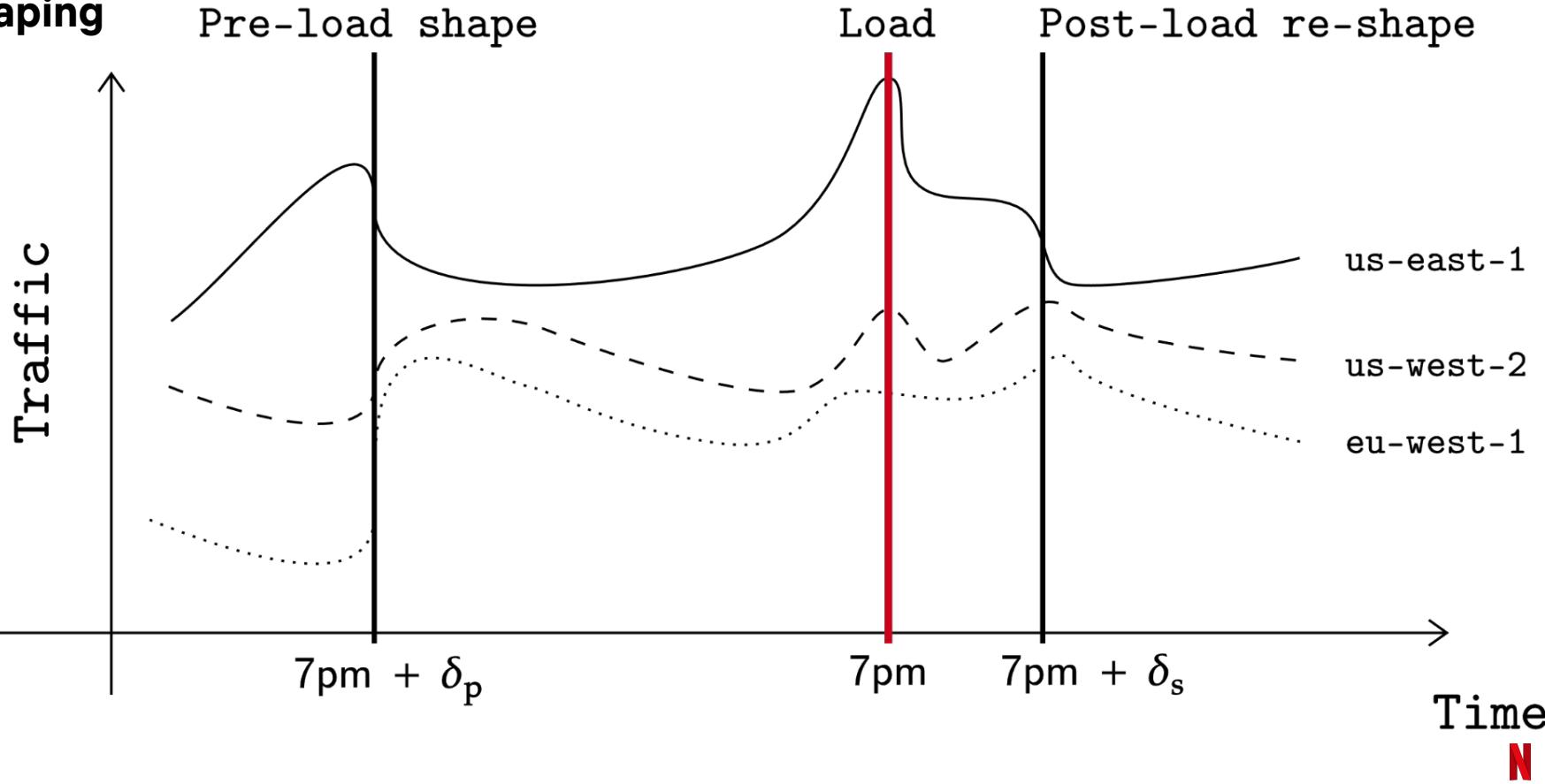
# SVOD usage is reasonably predictable



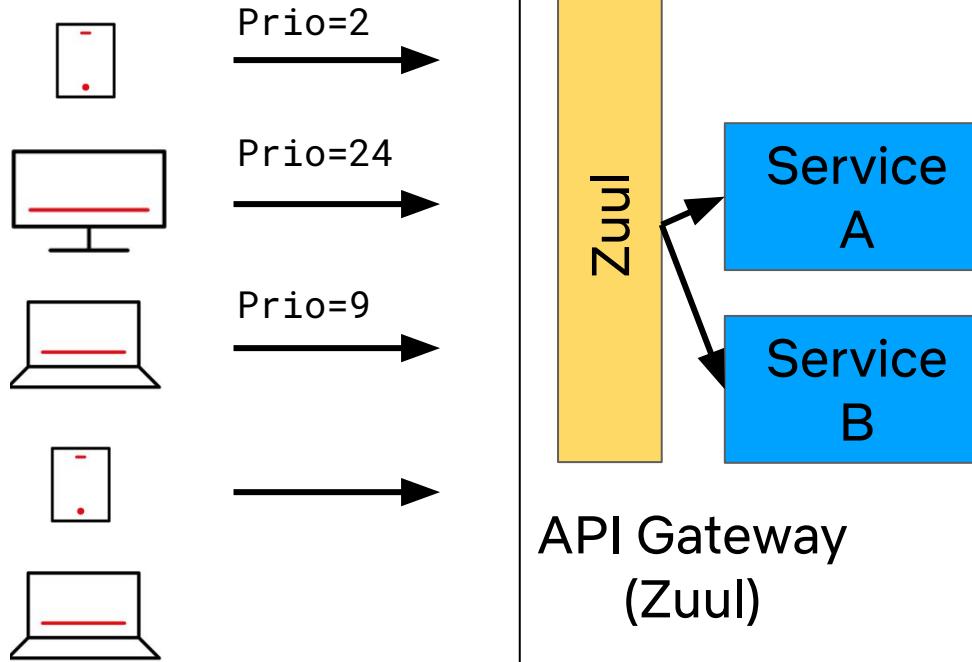
# Predictive Shaping



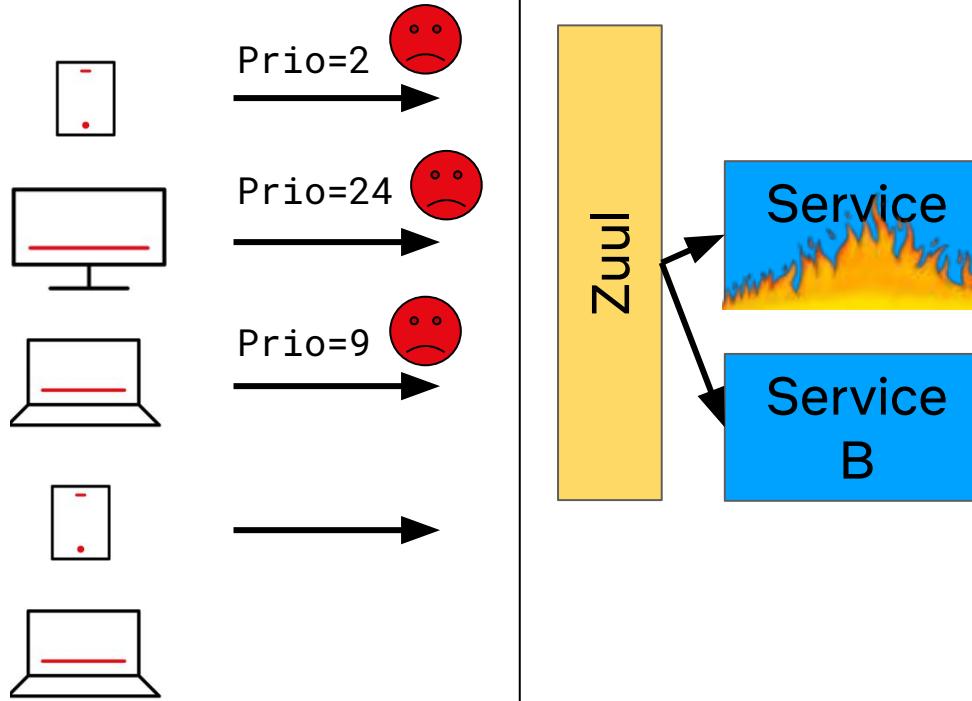
# Reactive Shaping



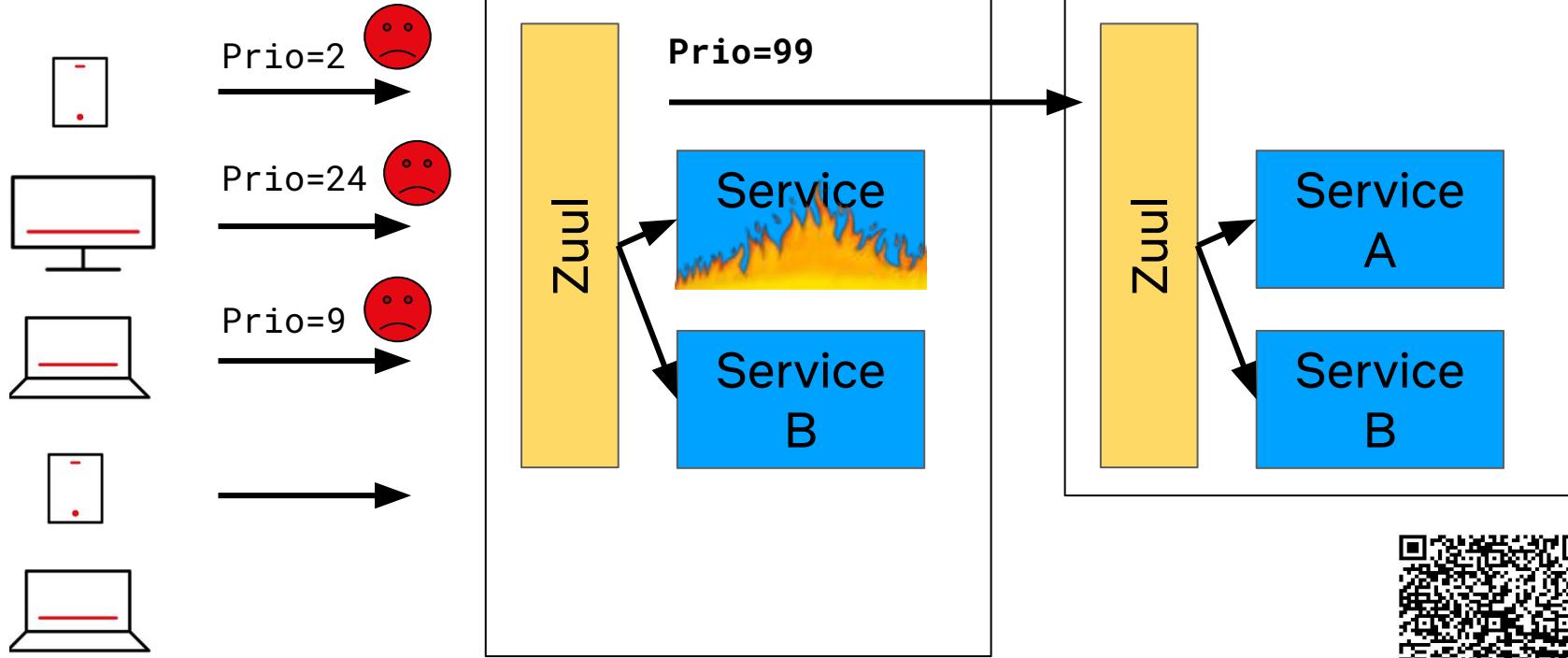
# Reactive Shifting



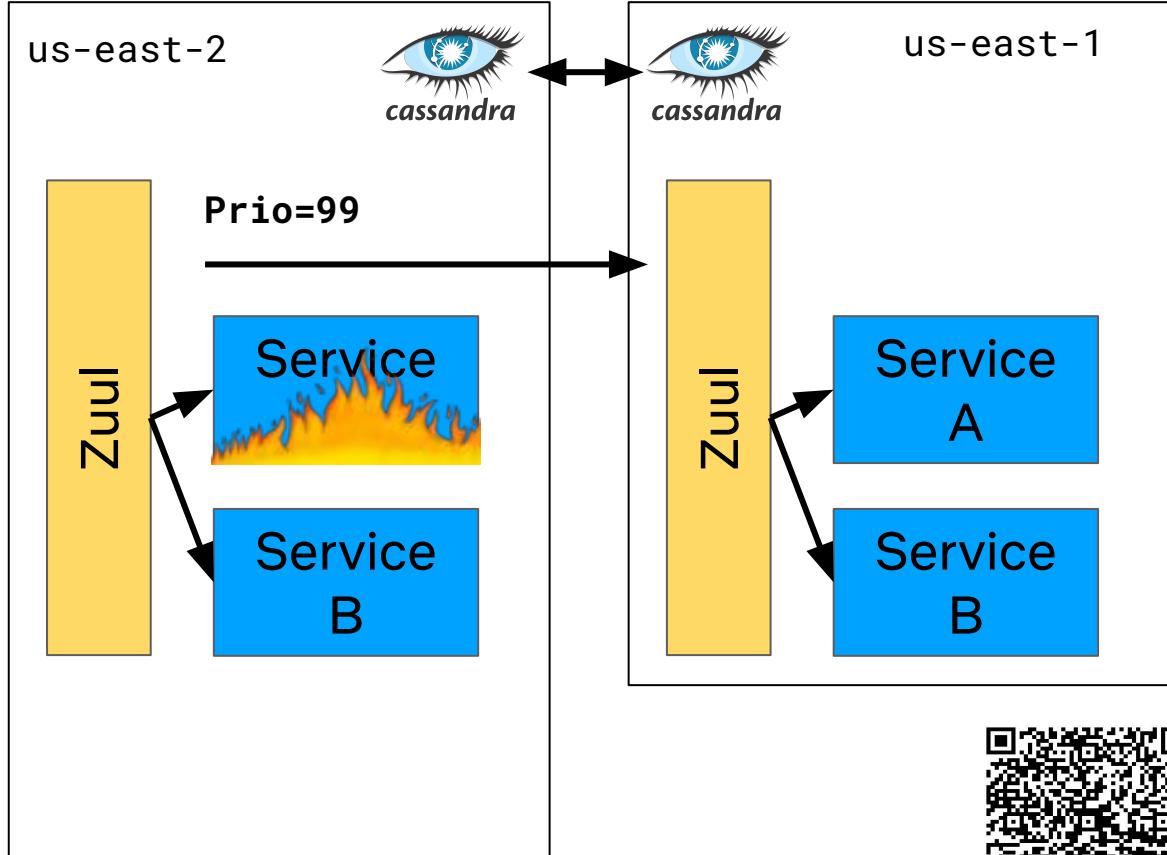
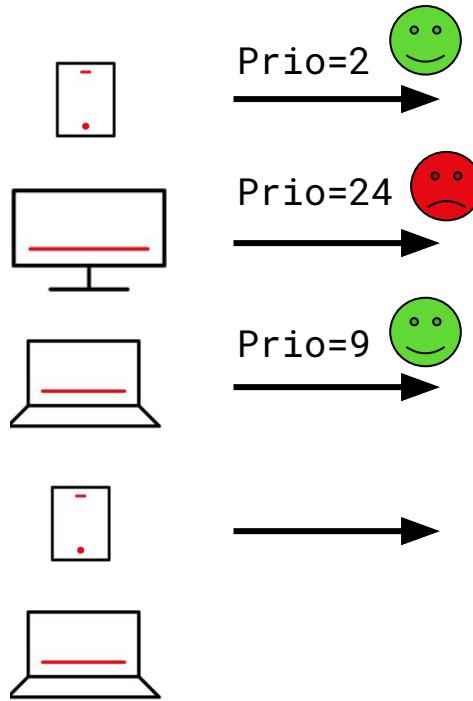
# Reactive Shifting



# Reactive Shifting

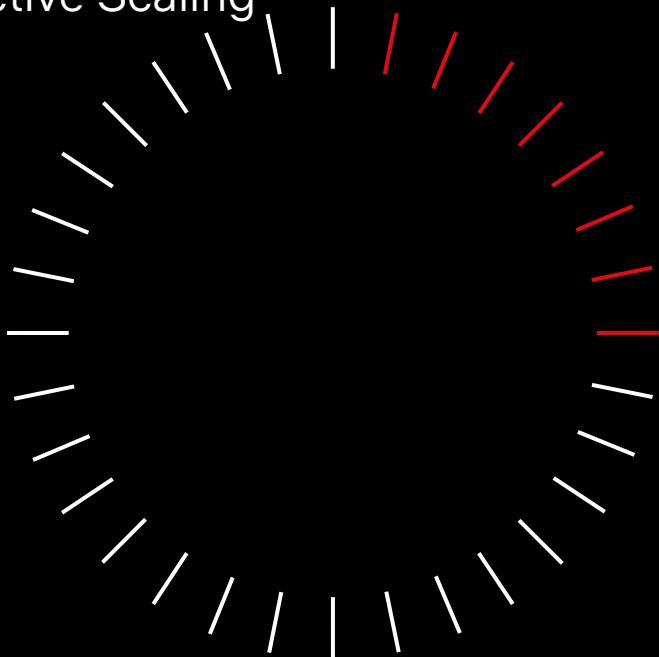


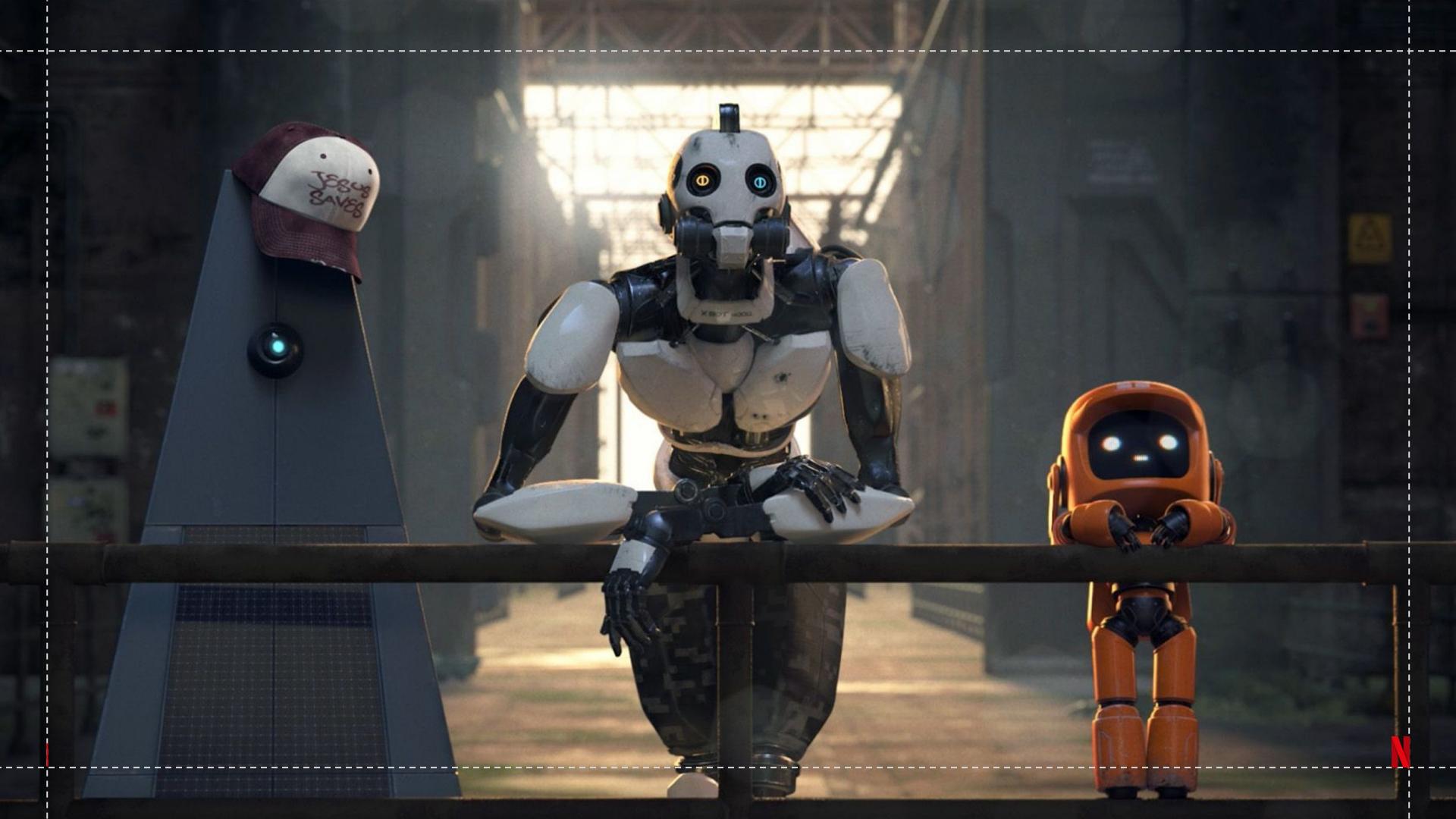
# Reactive Shifting



# Supply of Computers

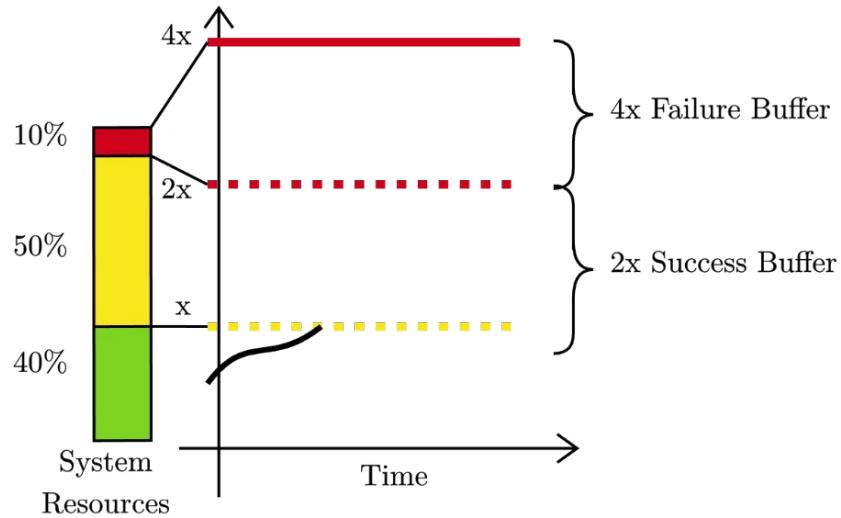
Cloud Realities  
Predictive Scaling  
Reactive Scaling



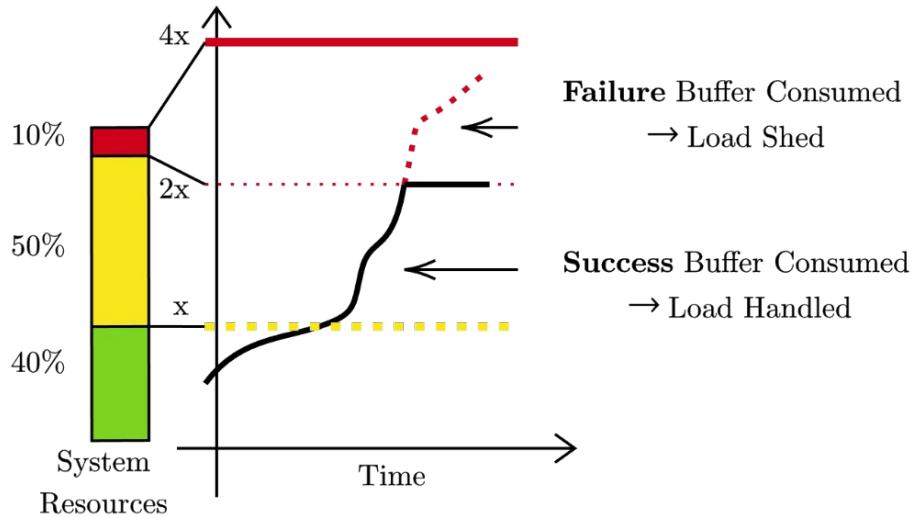


# Reason about headroom with Buffers

Normal System Load with Buffer

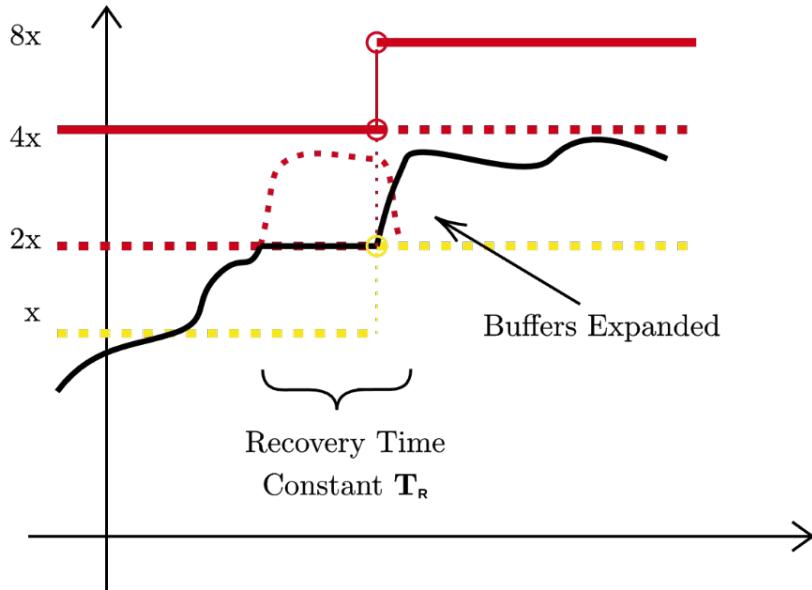


System Load Under Load Spike



# Buffer is Linked to Business Outcomes

Buffer Recovering After Load Spike

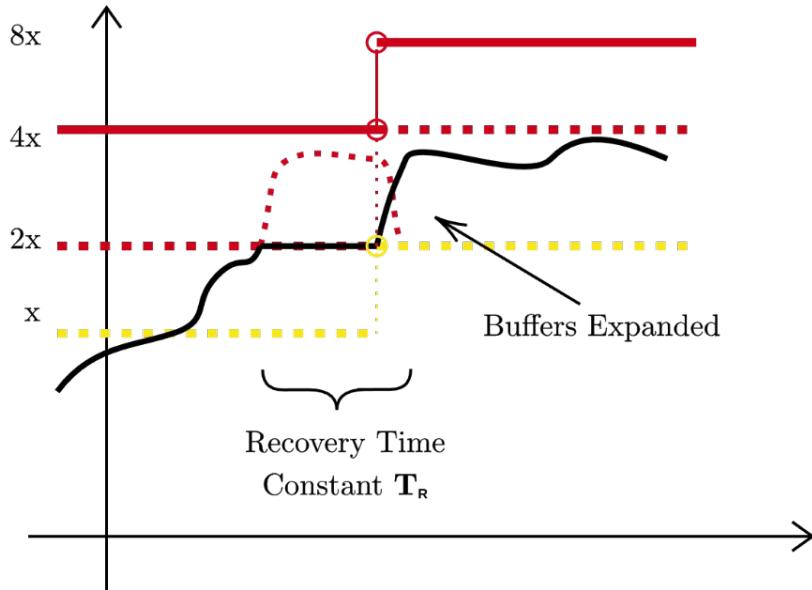


Buffer =  $F($   
Criticality Tier,  
Business Domain,  
Recovery Time Constant  
)

Low utilization is a tradeoff!

# Buffer is Linked to Business Outcomes

Buffer Recovering After Load Spike

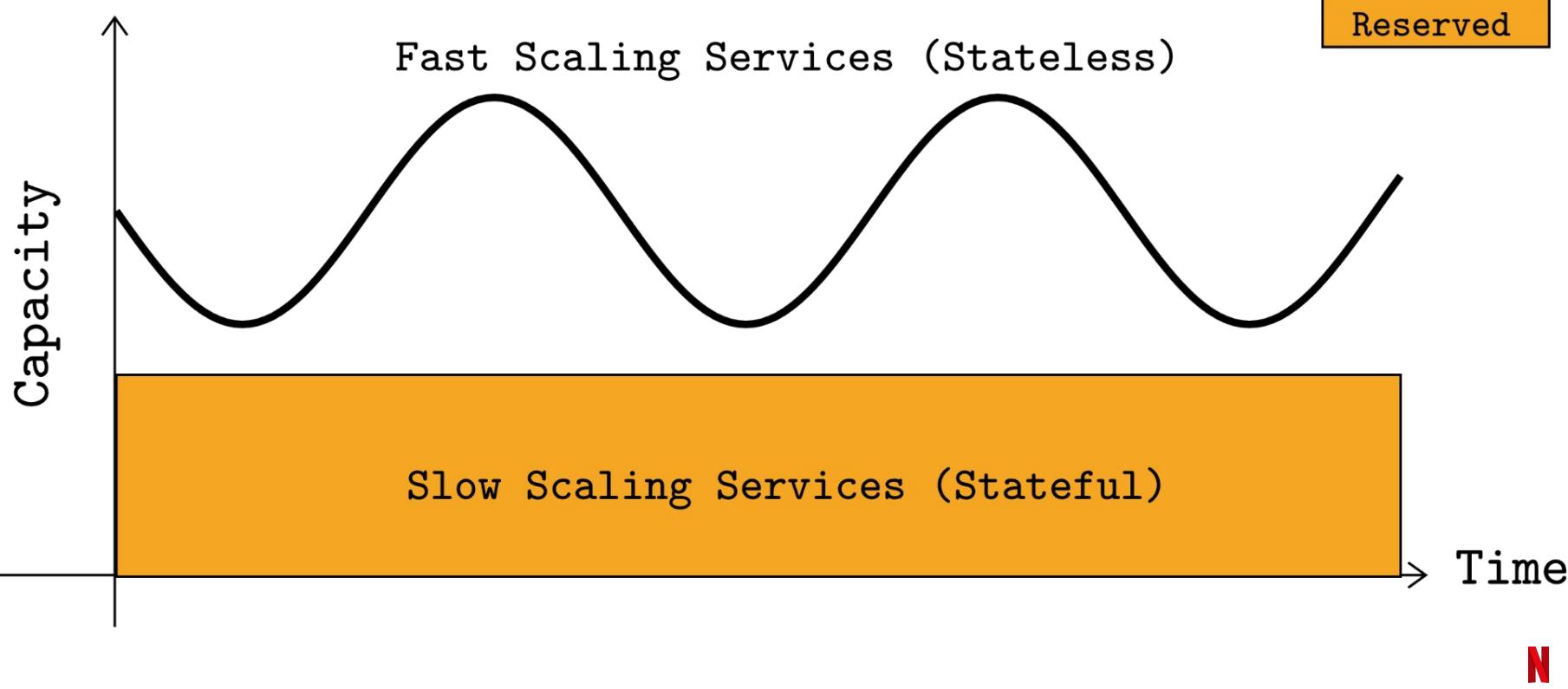


Buffer = F(  
Criticality Tier,  
Business Domain,  
Recovery Time Constant  
)

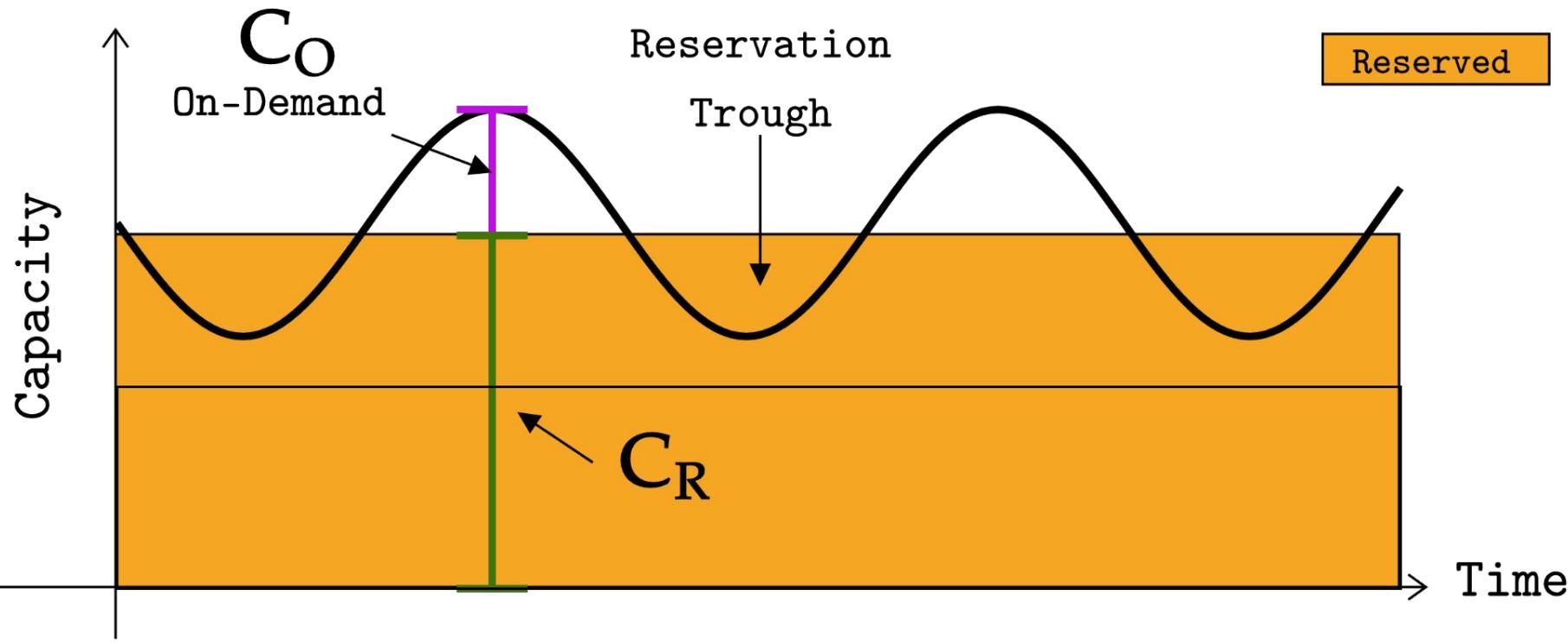
Low utilization is a tradeoff!

Fast recovering services need less buffer:  
**Stateless: 3-5m vs Stateful: 30-60m**

## Cloud Reality #1 Still have to Plan



## Cloud Reality #1 Still have to Plan



## Cloud Reality #1 Capacity Planning

It doesn't make financial sense to reserve everything

$$Capacity = C_R + C_O$$

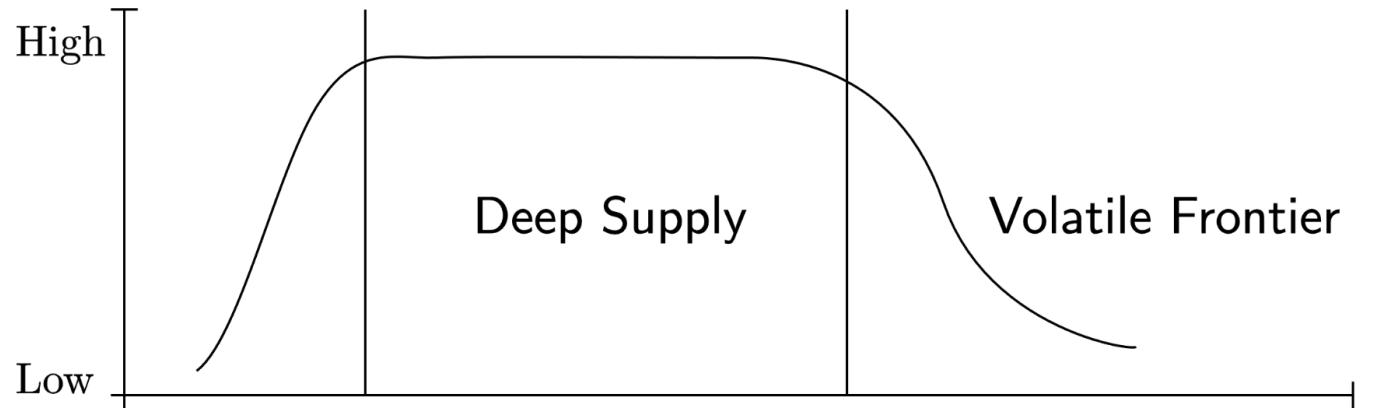
$$C_R \propto f(\text{time, pricing})$$

You will likely be using  
**On-Demand or Spot**

## Cloud Reality #2 Computers have *Variable Supply*

Have to carefully  
Adopt New Shapes

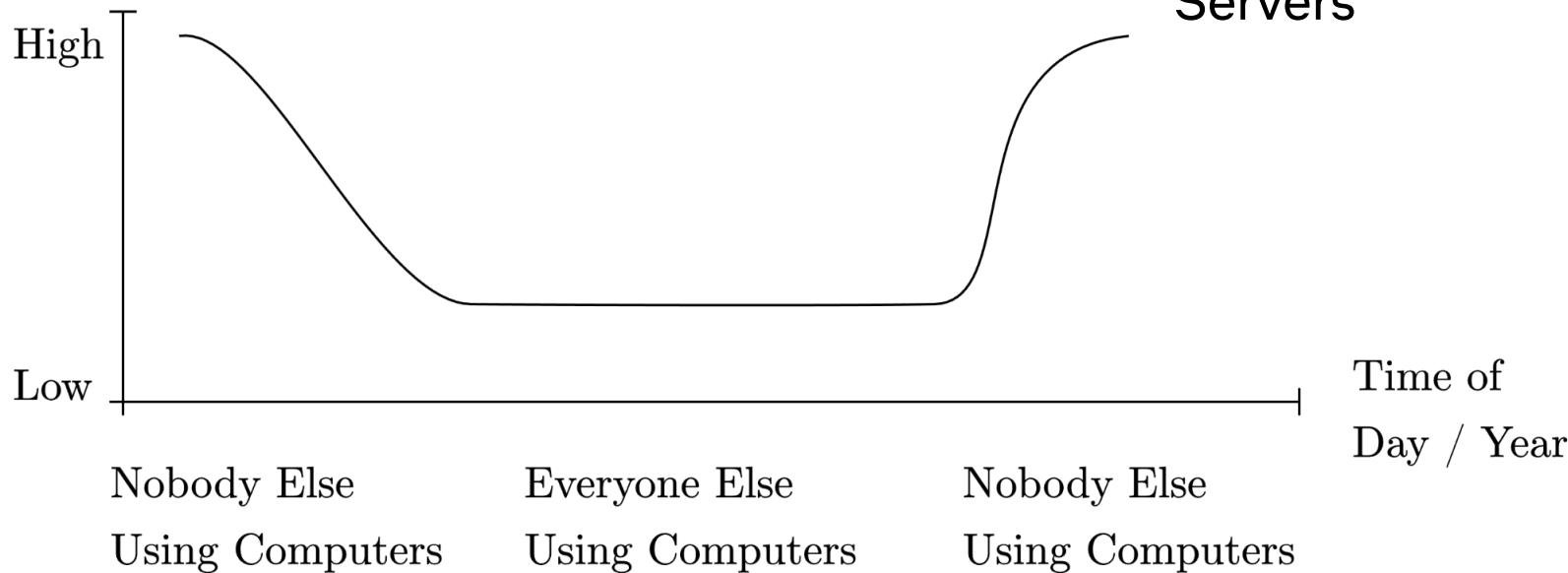
Compute Availability



	m4	i3	m5	i3en	m6i	i4i	m7i	m7a	i7ie	
Old	r4		r5		r6i		r7i	r7a		New
	i2		c5		c6i		c7i	c7a		

## Cloud Reality #2 Computers have *Variable Supply*

On-Demand Compute Availability



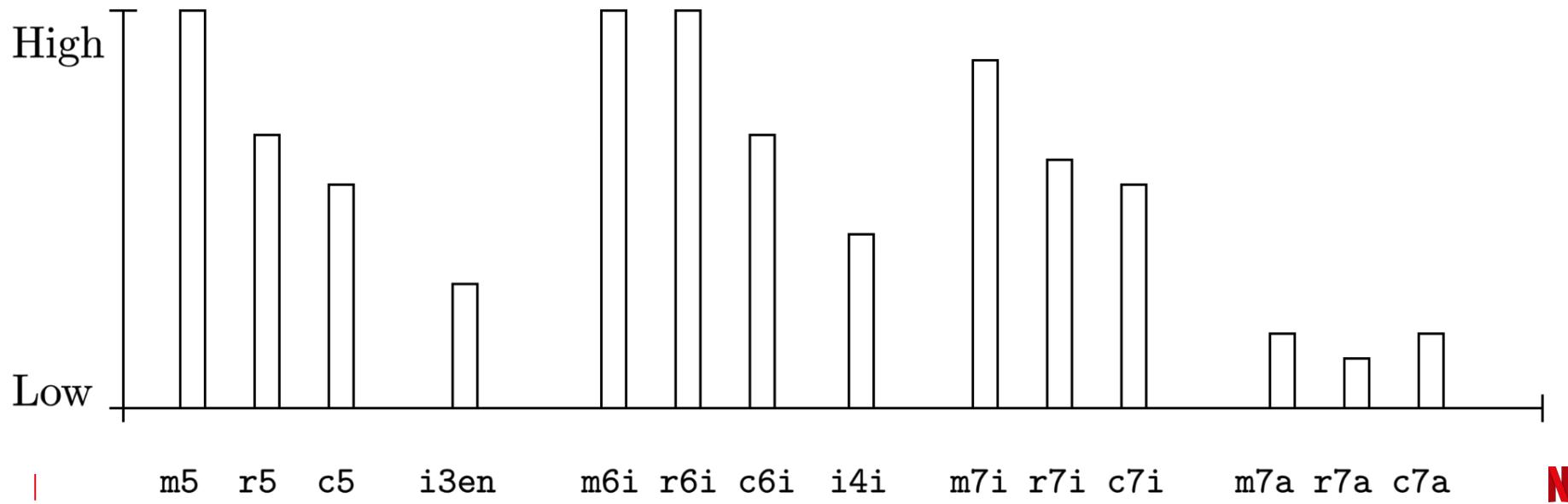
Options:

- Reservations
- On Demand Reservations
- Leverage Buffers
- Preferentially Allocating Servers

## Cloud Reality #2 Computers have *Variable Supply*

More flexibility == More Capacity  
Assign workloads by resources not names

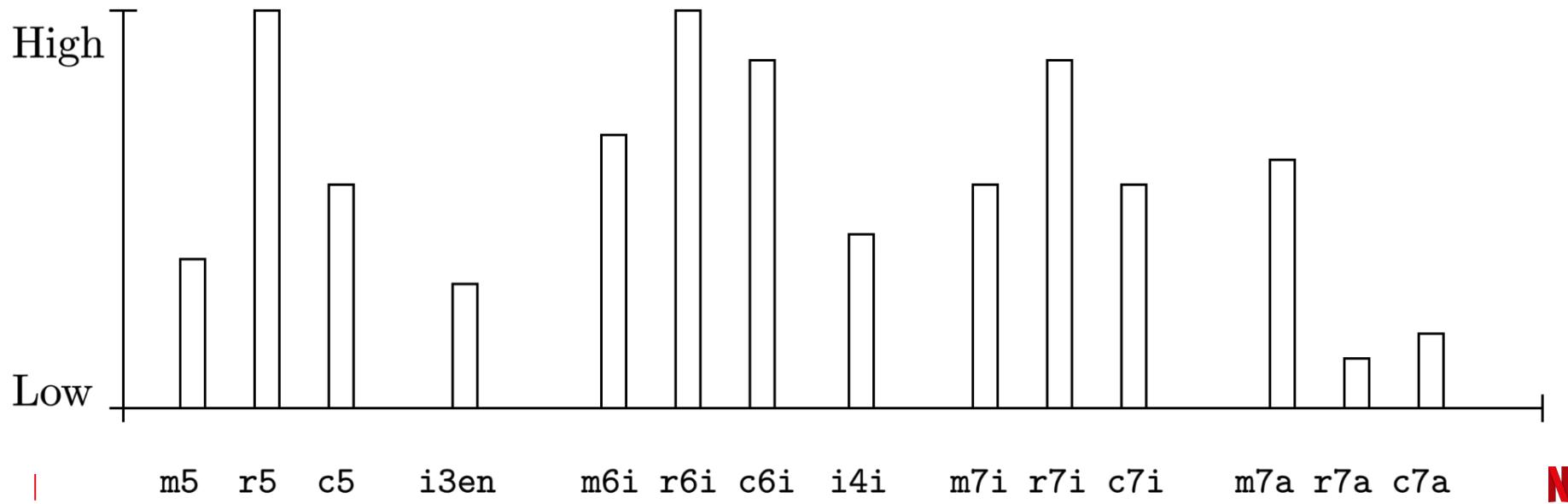
Compute Availability by Shape - Data Not Real



## Cloud Reality #2 Computers have *Variable Supply*

More flexibility == More Capacity  
Assign workloads by resources not names

Compute Availability by Shape - Data Not Real



## Cloud Reality #3 Different Servers are Meaningfully Different



```
"m7a.4xlarge": {  
    "name": "m7a.4xlarge",  
    "cpu": 16,  
    "cpu_cores": 16,  
    "cpu_ghz": 3.7,  
    "cpu_ipc_scale": 1.5,  
    "ram_gib": 61.04,  
    "net_mbps": 6250.0,  
    "drive": null  
},  
  
"m6id.4xlarge": {  
    "name": "m6id.4xlarge",  
    "cpu": 16,  
    "cpu_cores": 8,  
    "cpu_ghz": 3.5,  
    "cpu_ipc_scale": 1.0,  
    "ram_gib": 61.04,  
    "net_mbps": 6250.0,  
    "drive": {  
        "name": "ephem",  
        "size_gib": 885,  
        "read_io_per_s": 268332,  
        "write_io_per_s": 134168,  
        "single_tenant": false,  
        "read_io_latency_ms": {  
            "low": 0.1,  
            "mid": 0.125,  
            "high": 0.17,  
            "confidence": 0.9,  
            "minimum_value": 0.05,  
            "maximum_value": 2.0  
        }  
    }  
},
```

## Service Capacity Modeling

 Build passing

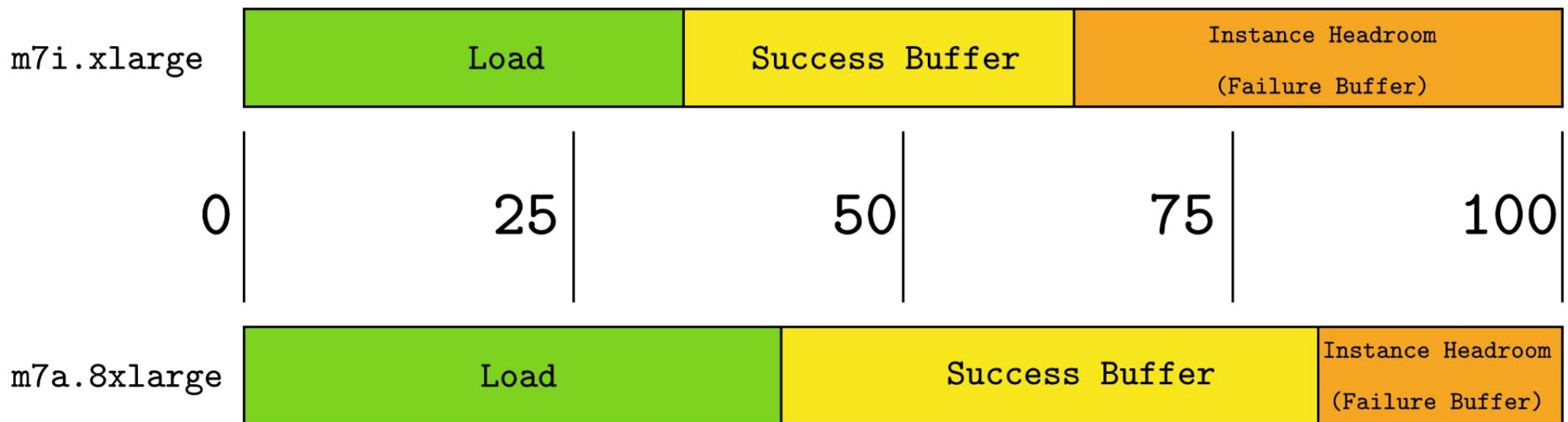
A generic toolkit for modeling capacity requirements in the cloud.

## Buffer Differs By Server Type

	Name	Failure Buffer %	Tier 1 Target %	Tier 0 Target %
19	m7a.xlarge	0.70	0.46	0.35
21	m7a.4xlarge	0.84	0.56	0.42
22	m7a.8xlarge	0.88	0.59	0.44
29	m7i.xlarge	0.62	0.41	0.31
31	m7i.4xlarge	0.79	0.53	0.40
32	m7i.8xlarge	0.85	0.57	0.42

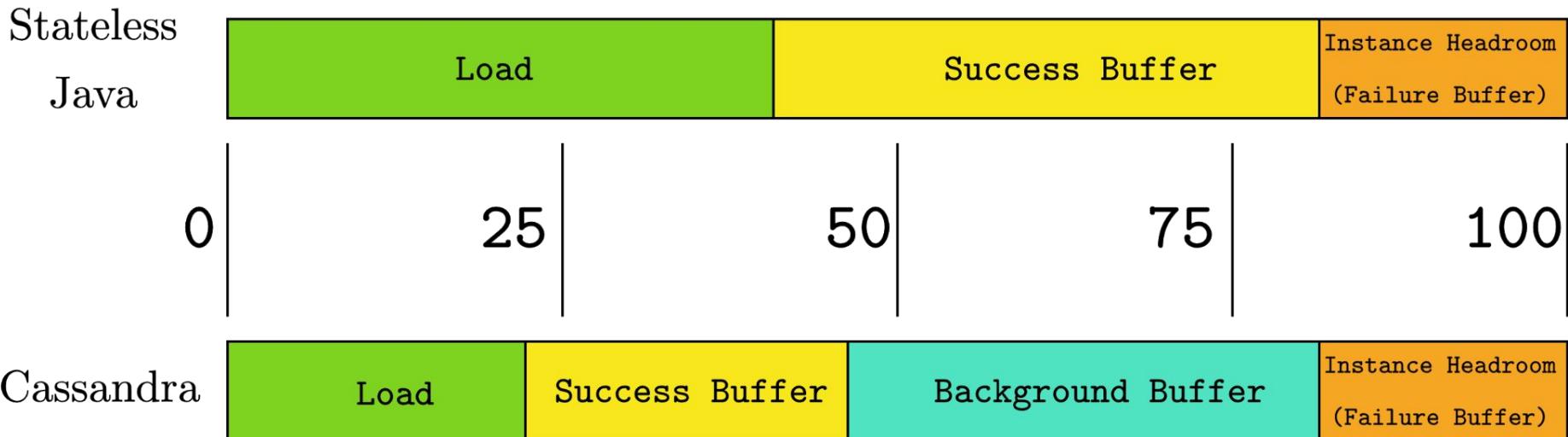
## Buffer Differs By Server Type

### Stateless Java CPU Utilization Buffers with 2x Success Buffer (Tier 0)



## Buffer Differs By Workload

### m7a.2x1 CPU Utilization Buffers with 2x Success Buffer



# Predict

Pre-scale services for demand  
Pre-scale stateful services

# React

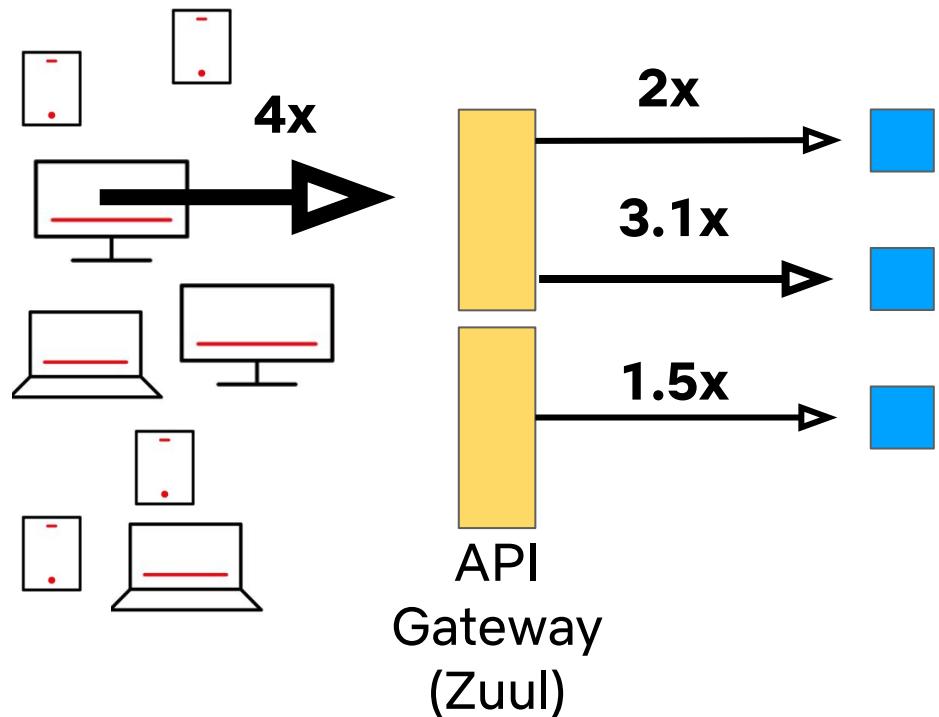
Autoscale out of trouble

## Pre-Scale for Load

Based on Projected Traffic, Pin Mins



**Different for  
Every Service**



Call Graph is not Uniform in Traffic or Criticality

**Tier 0 to Streaming Playback**

Tier 1 to Streaming Discovery

Tier 3 to Personalization

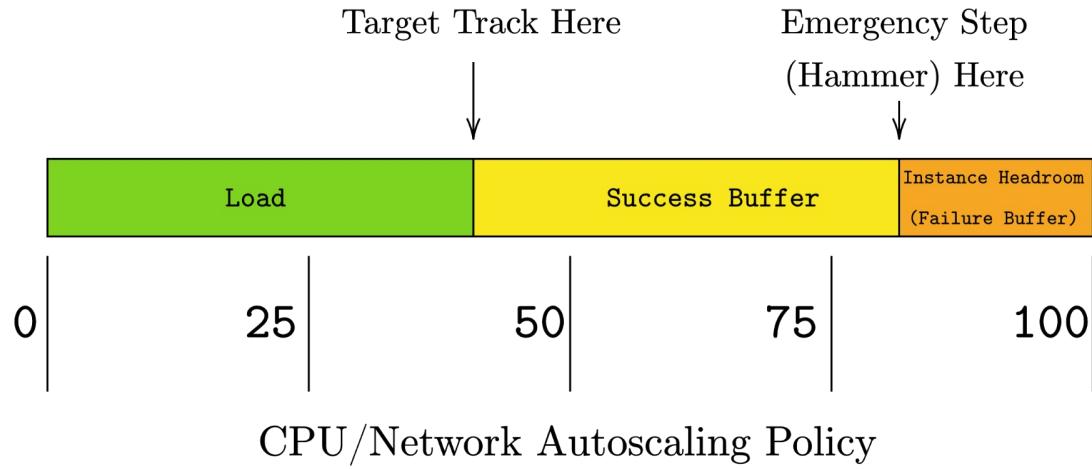
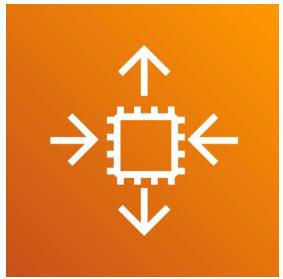


# What if we are wrong?

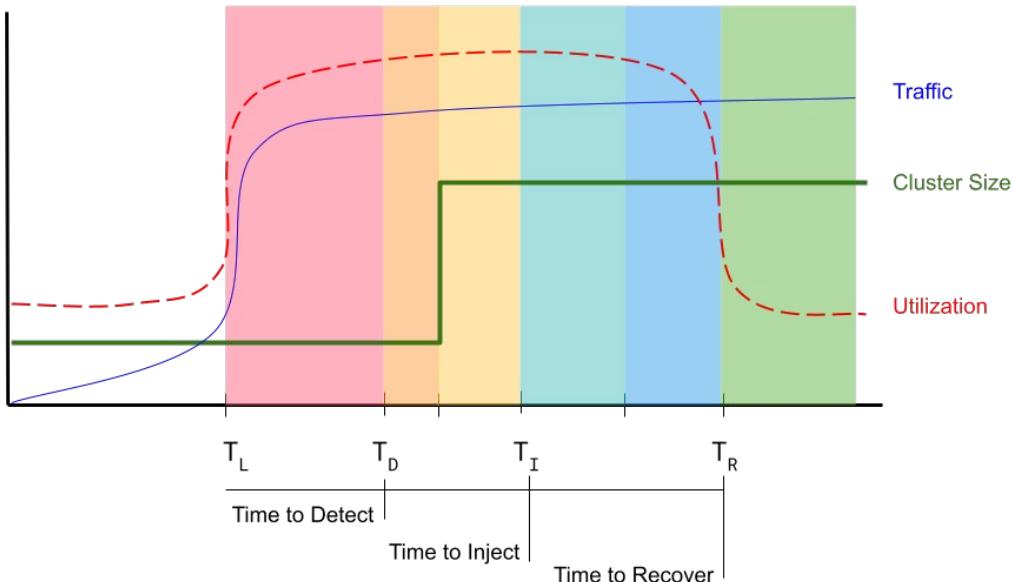


Perfectly Match  
Supply with Demand

# What if we are wrong?



# Autoscaling Can be Slow



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 OS 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  $T_R$ .

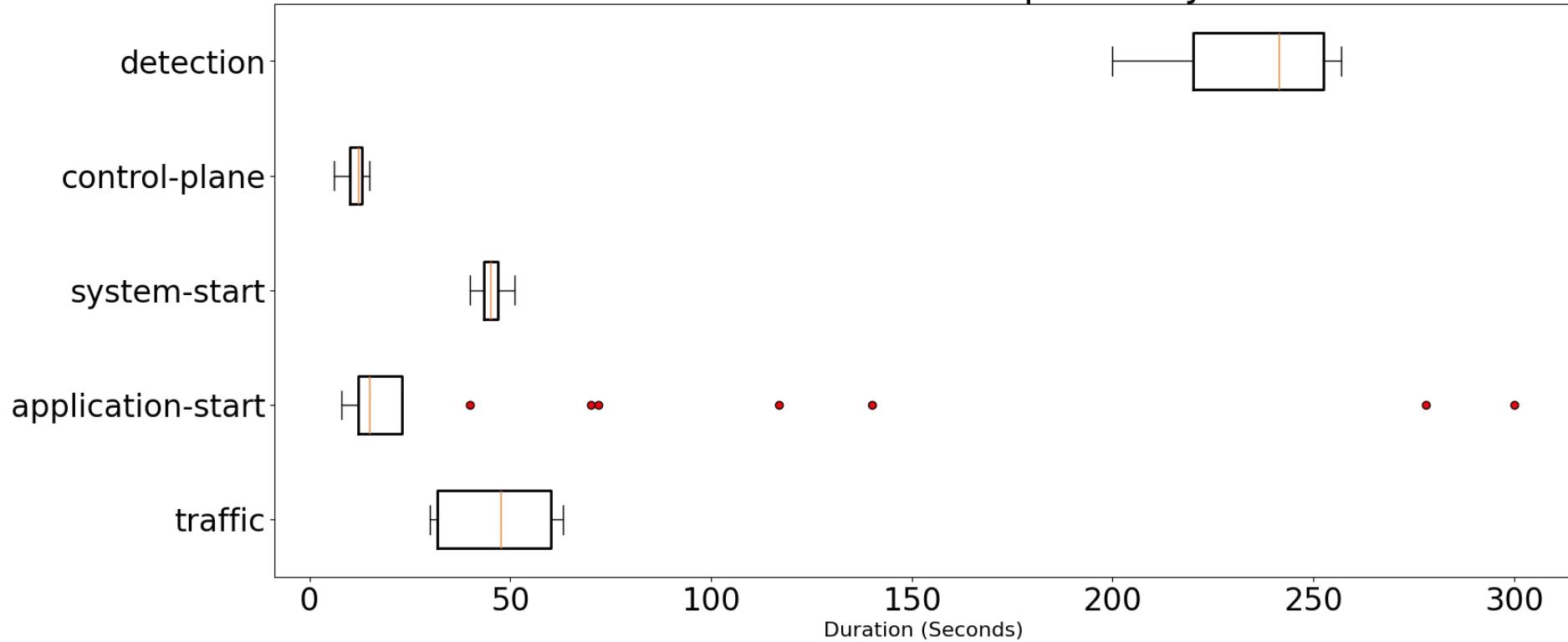
## Various Sources of Latency!

- Detection (Alert)
- Control Plane (Hardware)
- System Startup (Kernel)
- Application Startup (App)
- Traffic (Discovery)



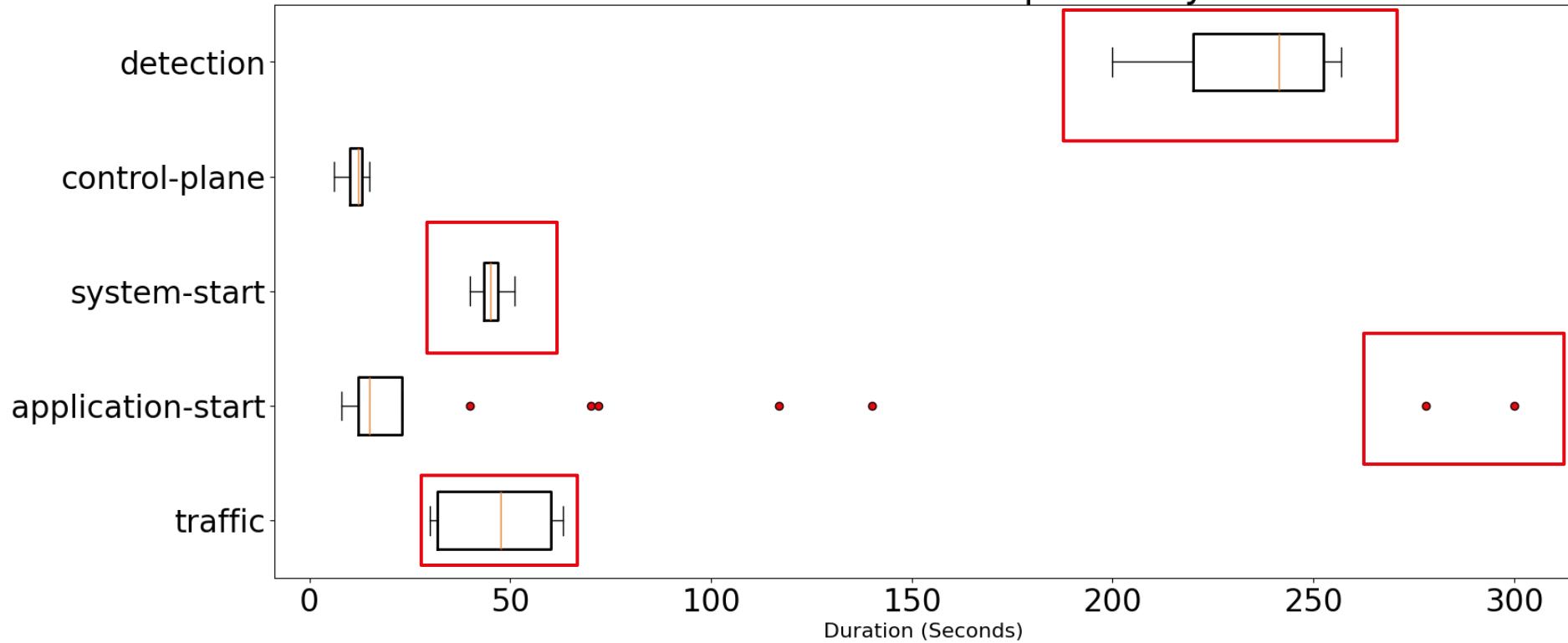
# Break it Down

Breakdown of Startup Latency



# Break it Down

Breakdown of Startup Latency



# Solve Each Part



High Resolution Metrics

Observe actual start latency (cooldown)

Observe RPS/CPU and Hammer on RPS ladder.

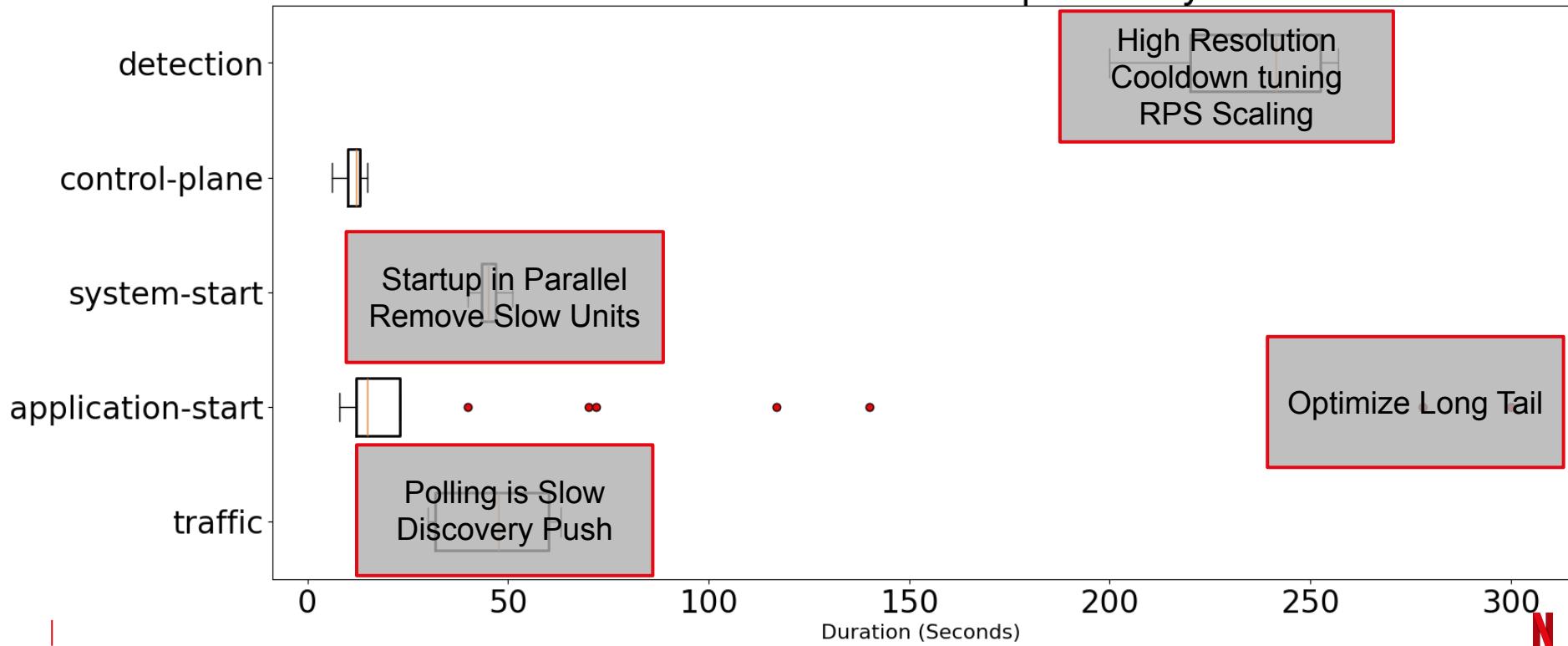
detection





## Break it Down

### Breakdown of Startup Latency



# Does This Really Work?



**Tier 0 Before Tuning  
10x Load Spike TTR**

8-15M

Snail: <https://commons.wikimedia.org/wiki/File:Snail.jpg>

# Does This Really Work?



~70% reduction!



**Tier 0 Before Tuning**  
**10x Load Spike TTR**

8-15M

**Tier 0 After Tuning**  
**10x Load Spike TTR**

3-4M

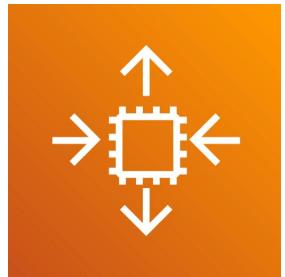
# Stateless Resilience

Load Shedding  
CPU  
IO  
Prioritization

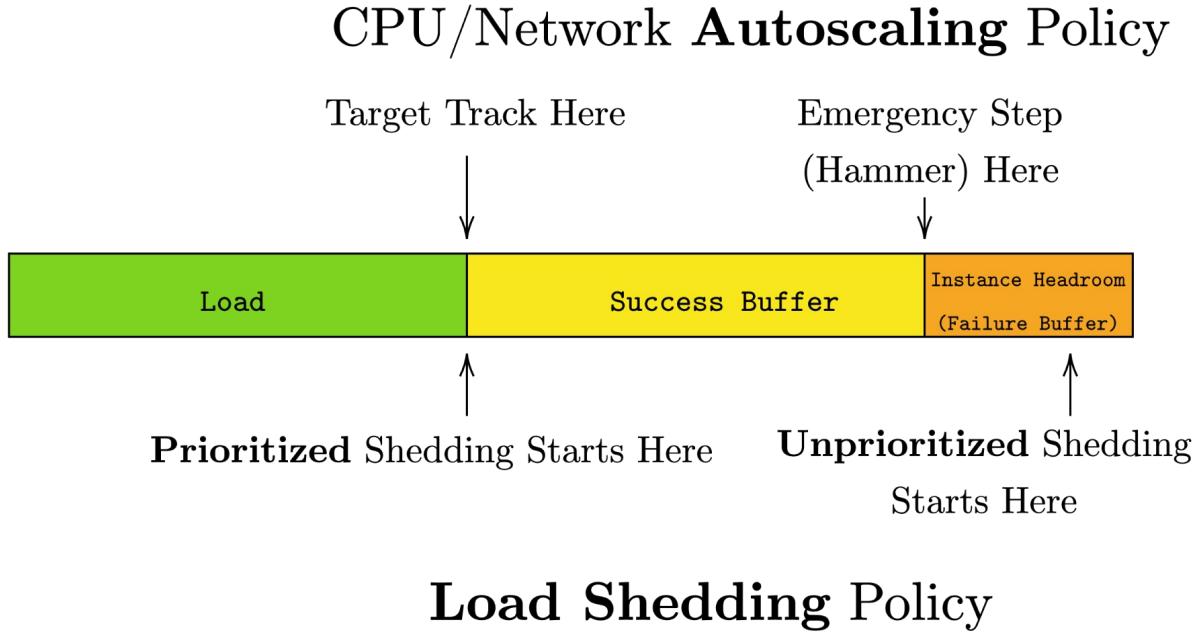
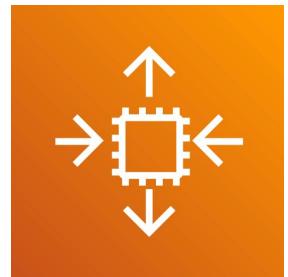


N

# What to do *While* we are Wrong?

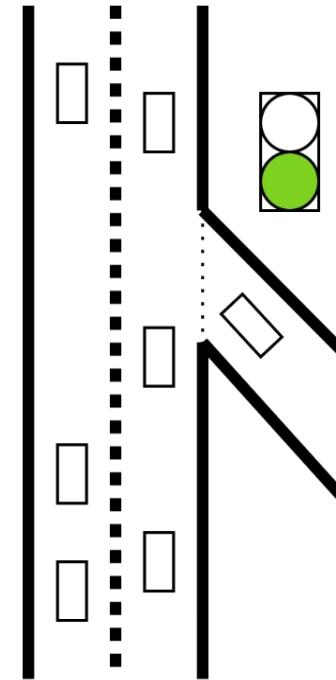
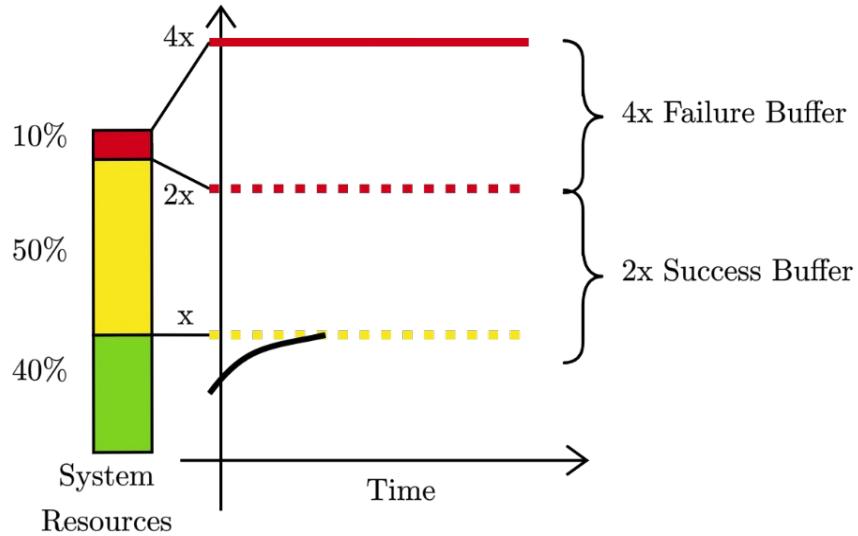


# What to do *While* we are Wrong?



# CPU Load Sheding

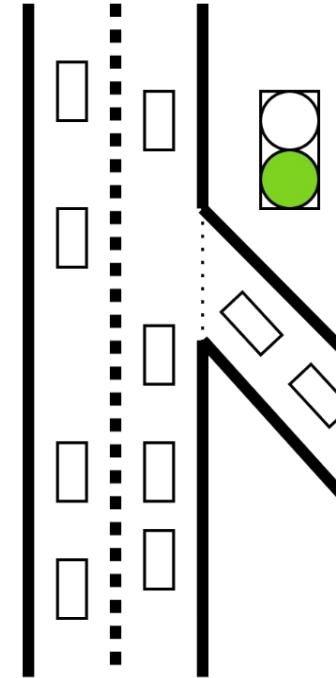
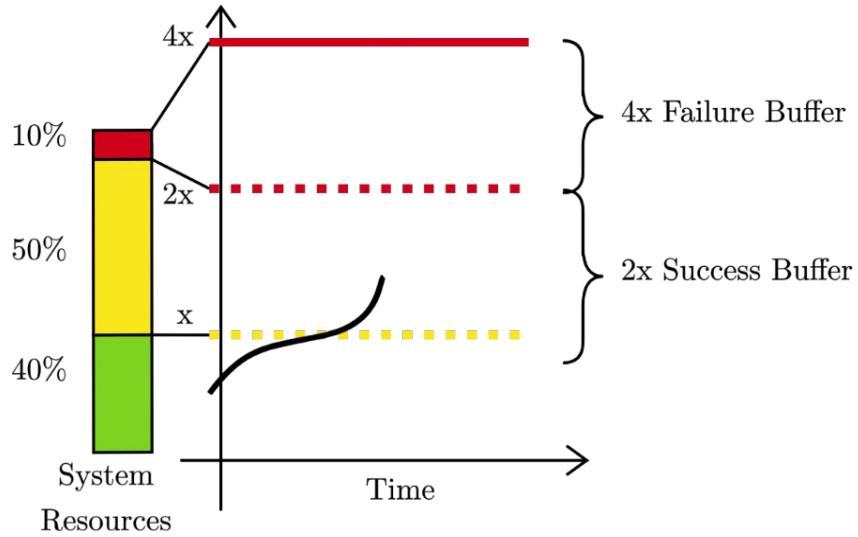
Normal System Load with Buffer



$$T_{transit} = 10m$$

# CPU Load Sheding

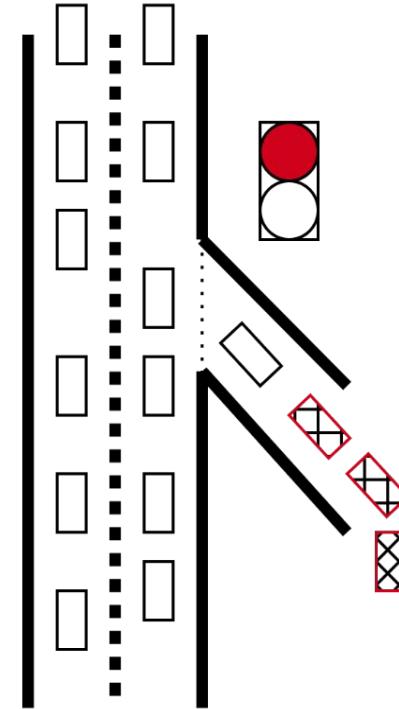
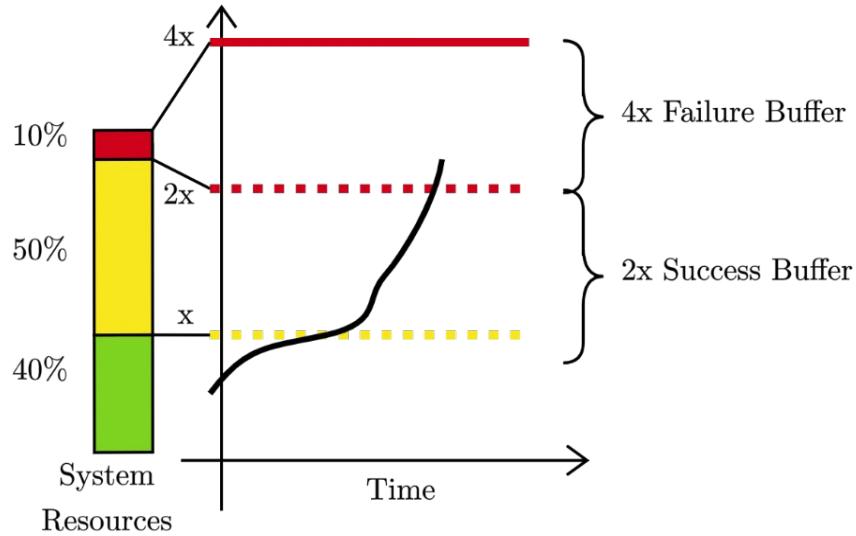
Normal System Load with Buffer



$$T_{transit} = 15m$$

# CPU Load Sheding

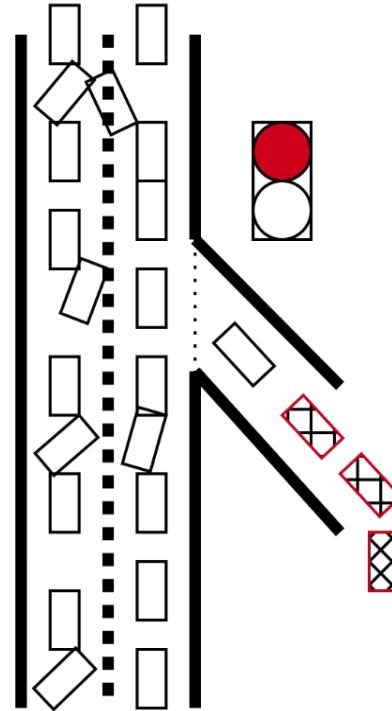
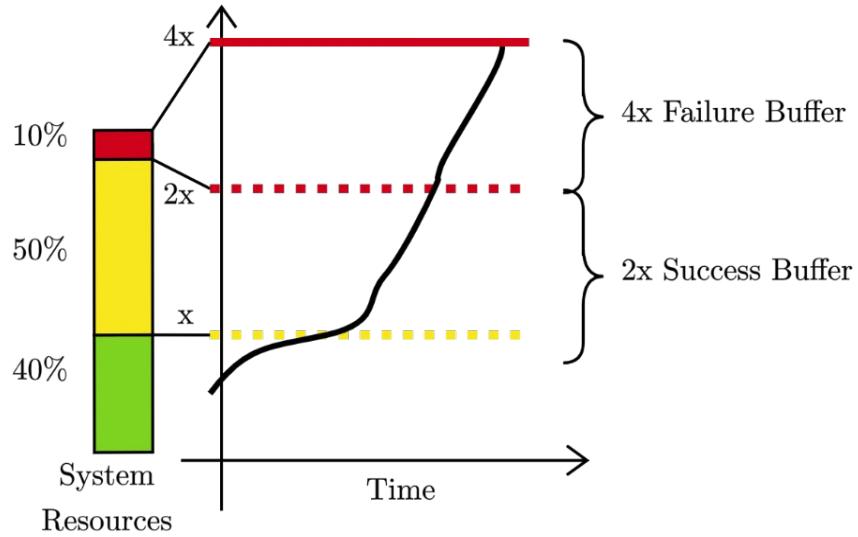
Normal System Load with Buffer



$$T_{transit} = 30m$$

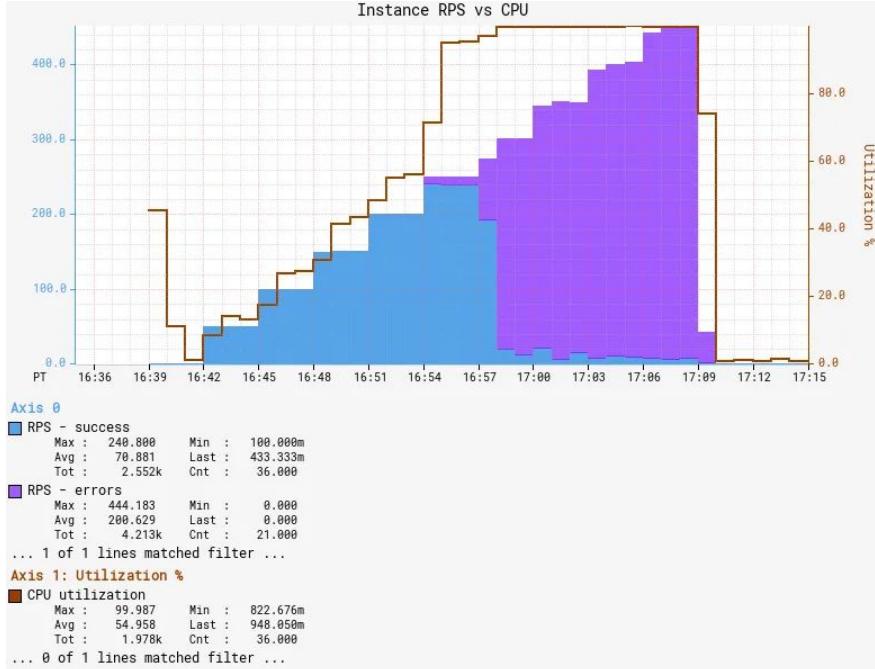
# CPU Load Sheding

Normal System Load with Buffer



$$T_{transit} = \infty$$

# Congestive Failure is Bad



Congestive Failure, a.k.a

"Queueing Without Bound"

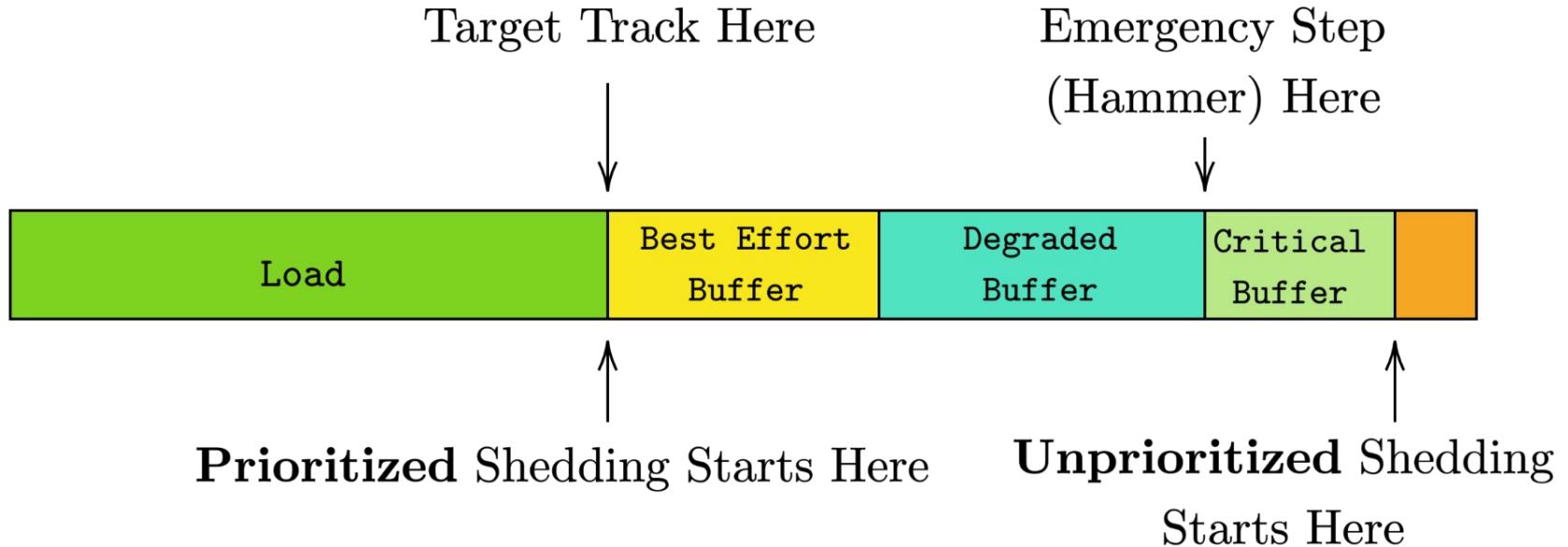
"Falling Over"

"Total System Failure"

"A real bad time™"

## Prioritize Your Success Buffer

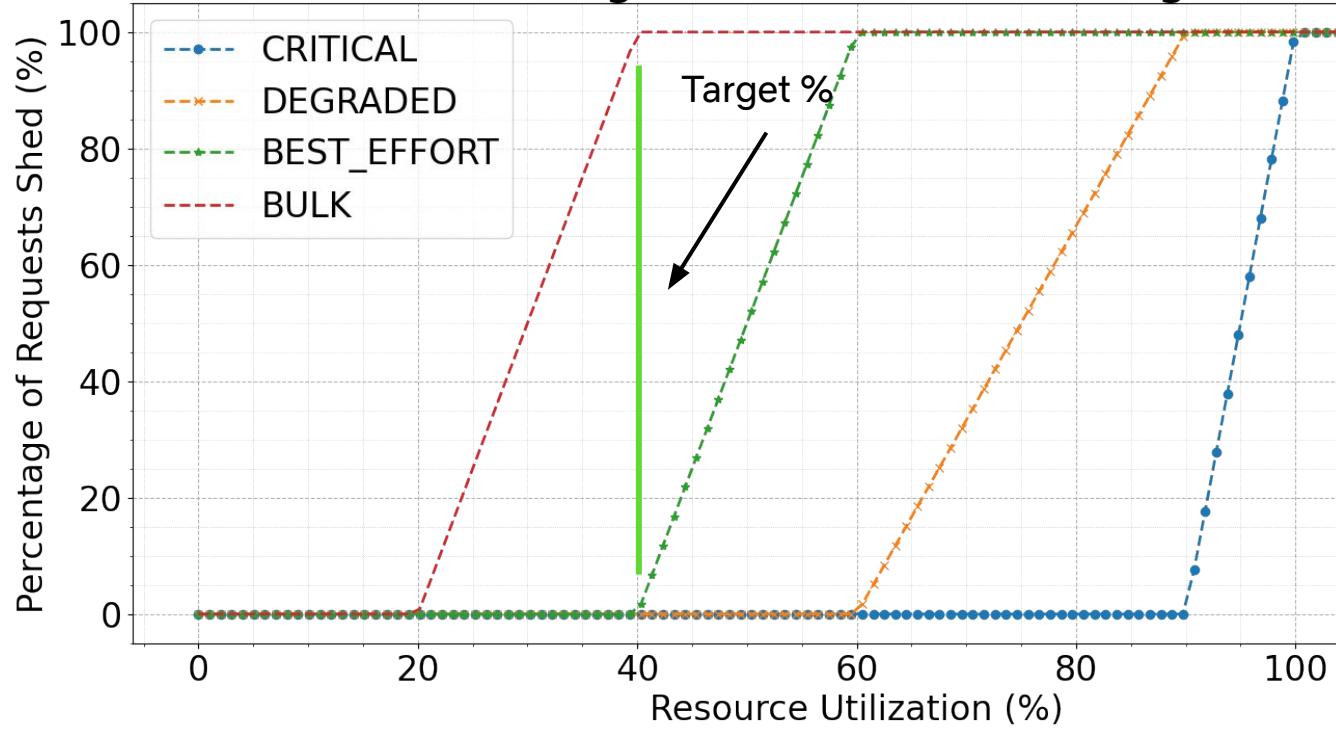
# CPU/Network Autoscaling Policy



Prioritized Load Shedding Policy

# Prioritization Creates More Buffer

## Progressive Load Shedding



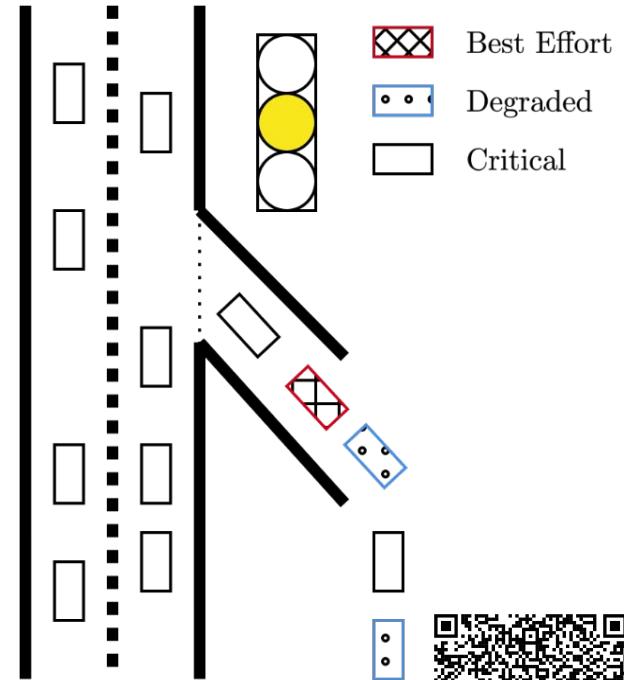
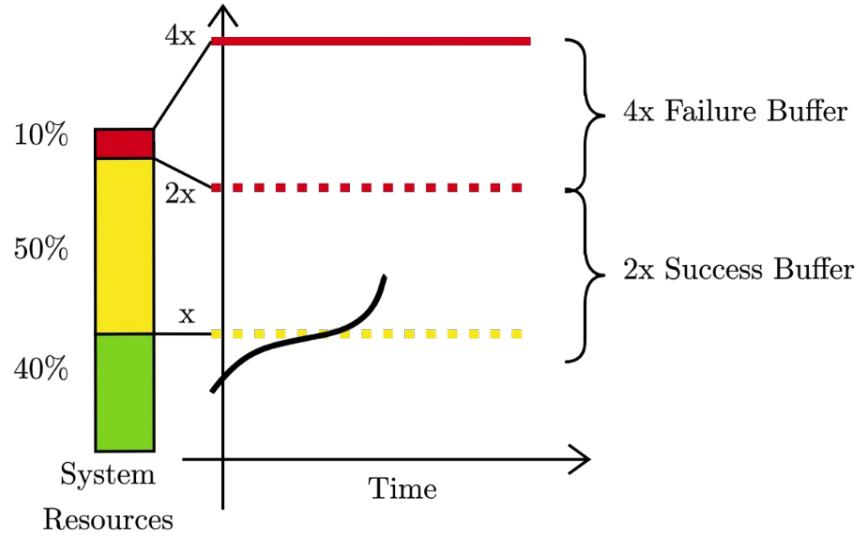
# Prioritization Creates More Buffer

```
return context -> {
    Request req = context.getRequest();
    // Creates More
    // 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());
}
```



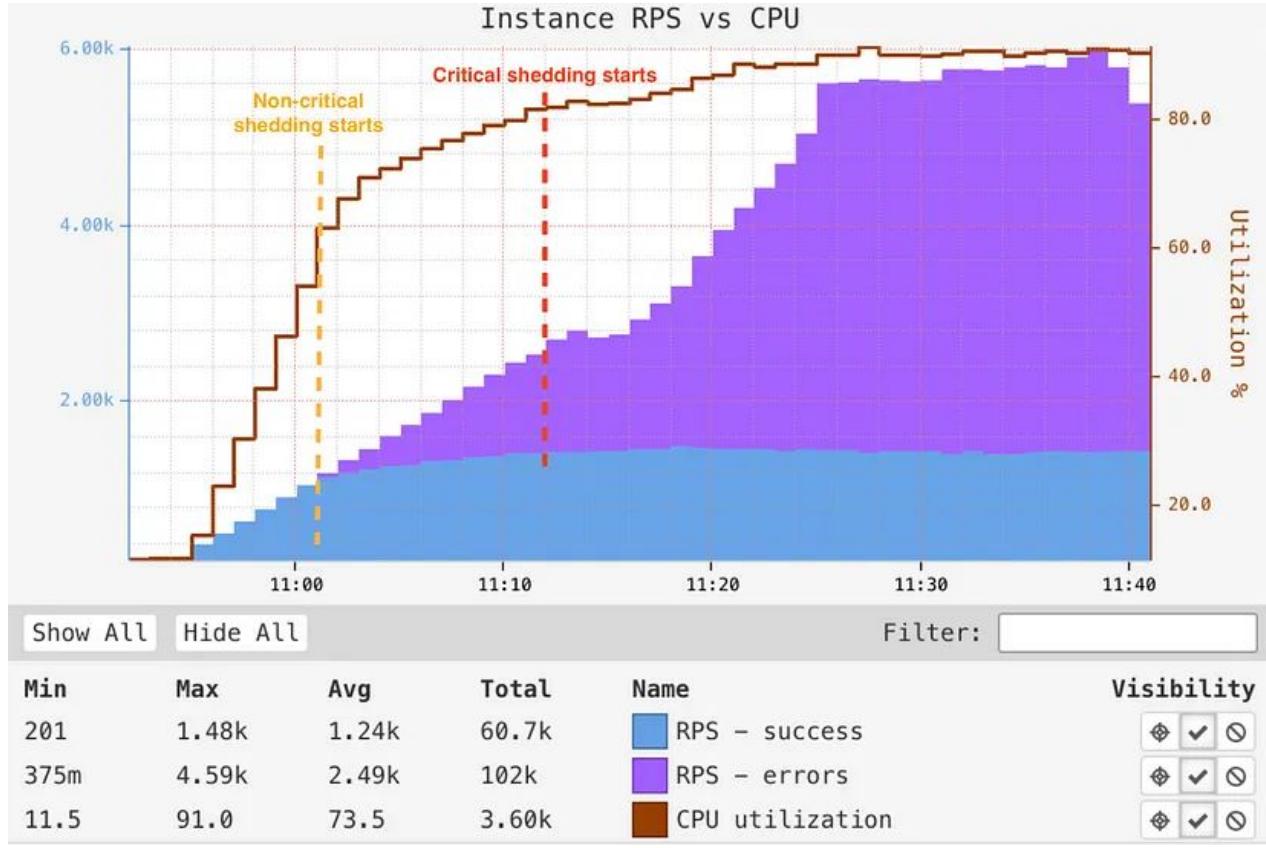
# Prioritization Creates More Buffer

Normal System Load with Buffer



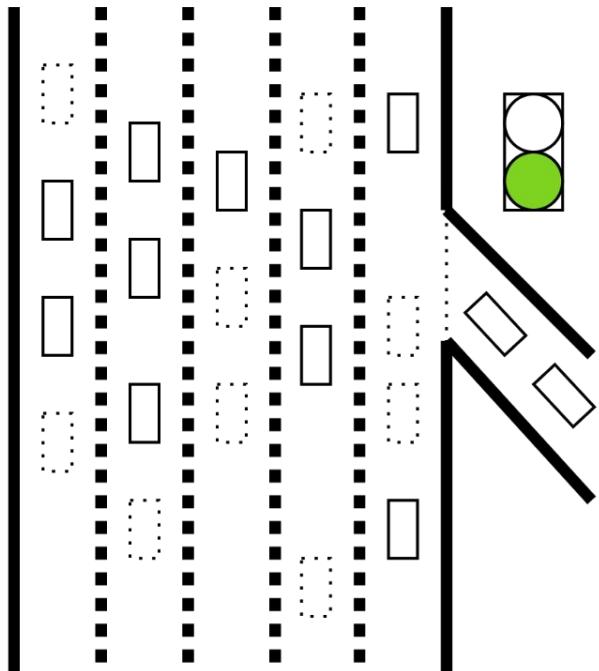
$$T_{transit} = 15m$$

# Prioritization Creates More Buffer



## IO Load Shedding

# How do we shed for IO workloads?



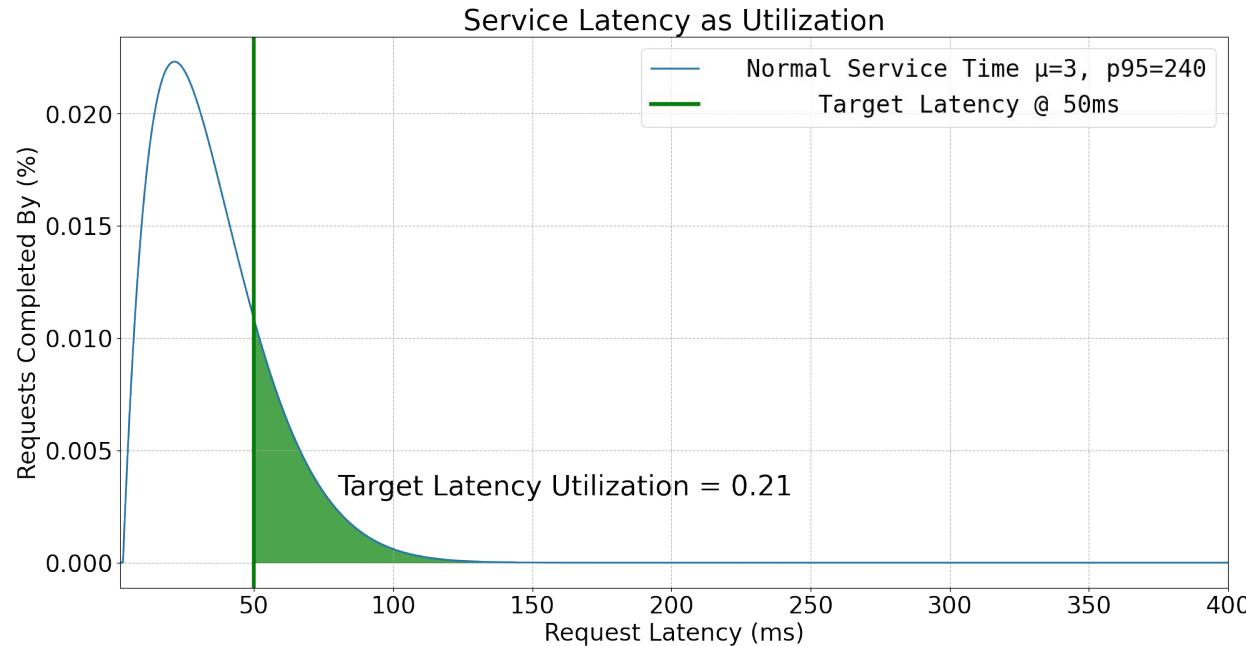
$$T_{transit} = 15m$$

$$T_{transit} = 45m$$

## IO Load Shedding

```
utilization(namespace) = {  
    overall = 21  
    latency = {  
        slo_target = 21,  
        slo_max = 0  
    }  
    system = {  
        storage = 17,  
        compute = 10,  
    }  
}
```

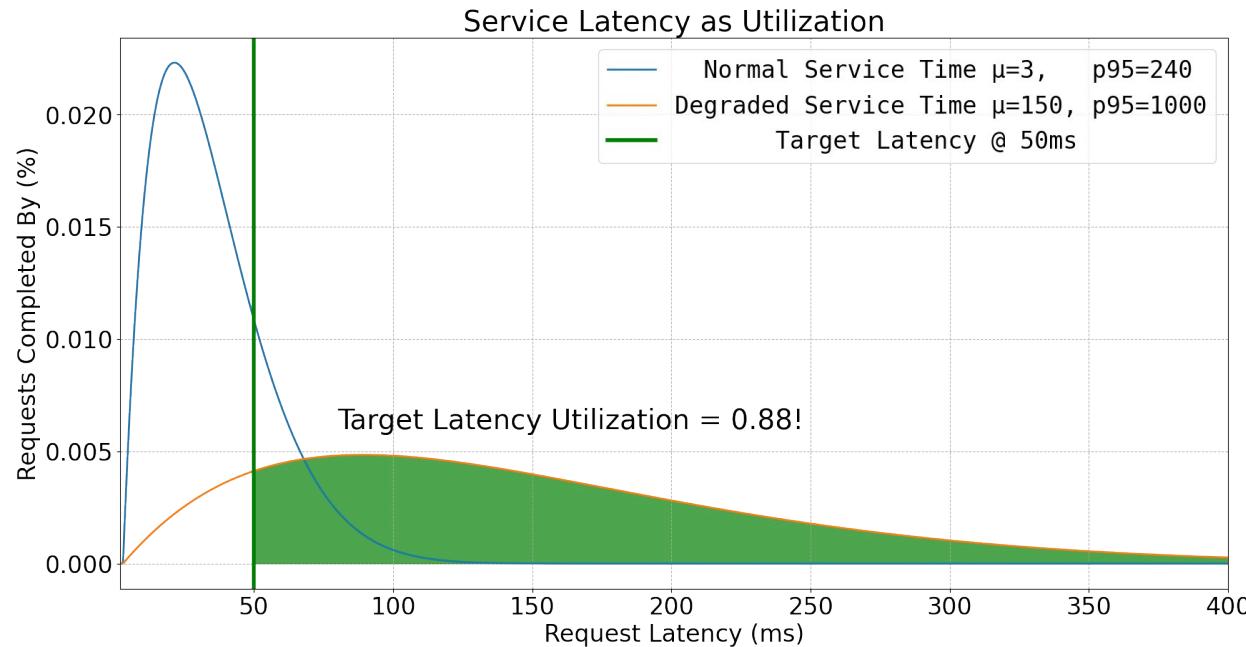
# Use Latency Service-Level-Objective Utilization as a Proxy!



# IO Load Shedding

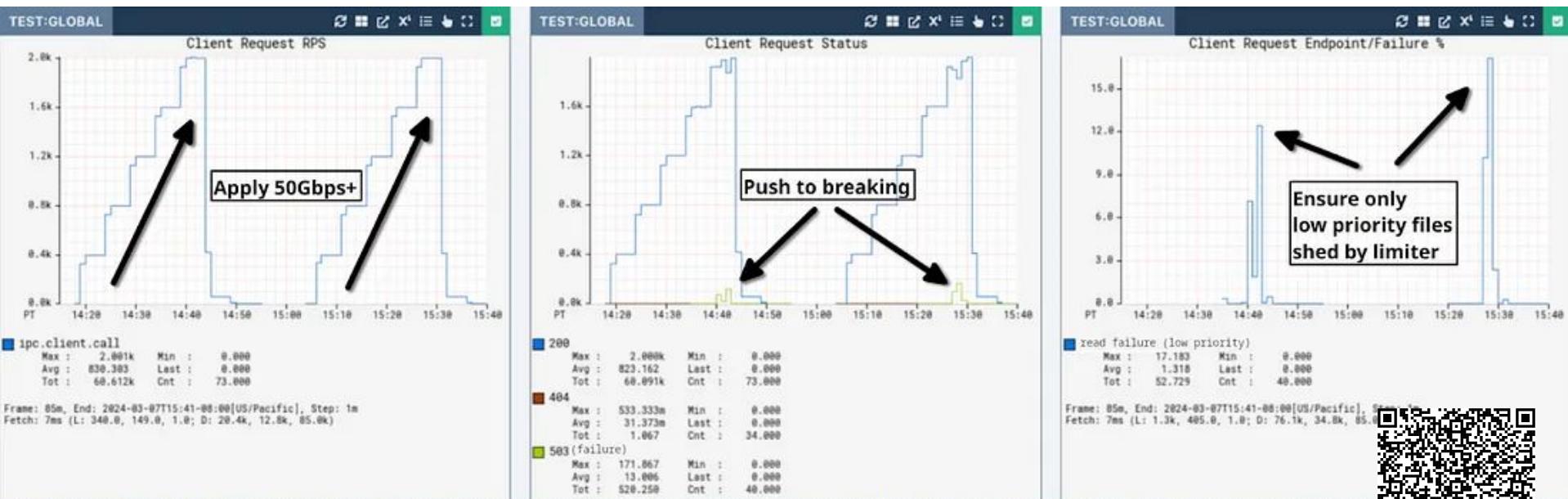
```
utilization(namespace) = {  
    overall = 88  
    latency = {  
        slo_target = 88,  
        slo_max = 5  
    }  
    system = {  
        storage = 17,  
        compute = 10,  
    }  
}
```

## Use Latency Service-Level-Objective Utilization as a Proxy!

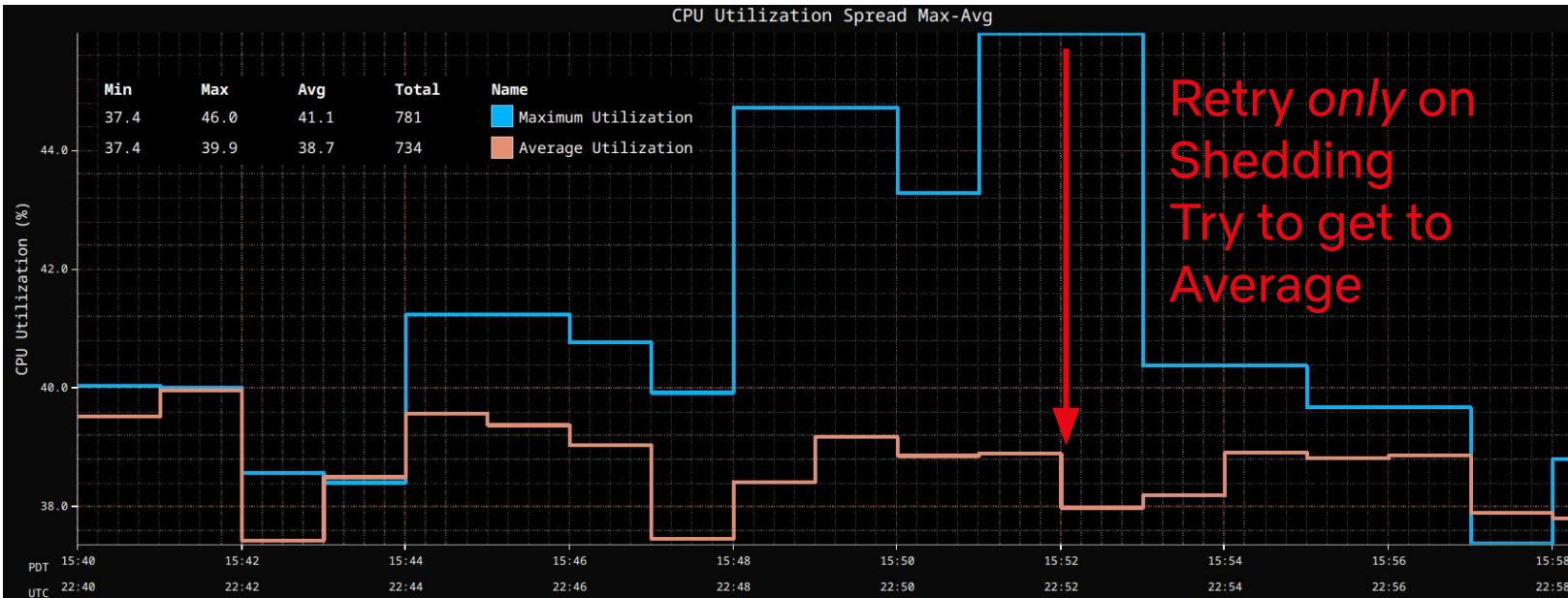


# IO Load Shedding

## Use Latency Service-Level-Objective Utilization as a Proxy!



# Retries?



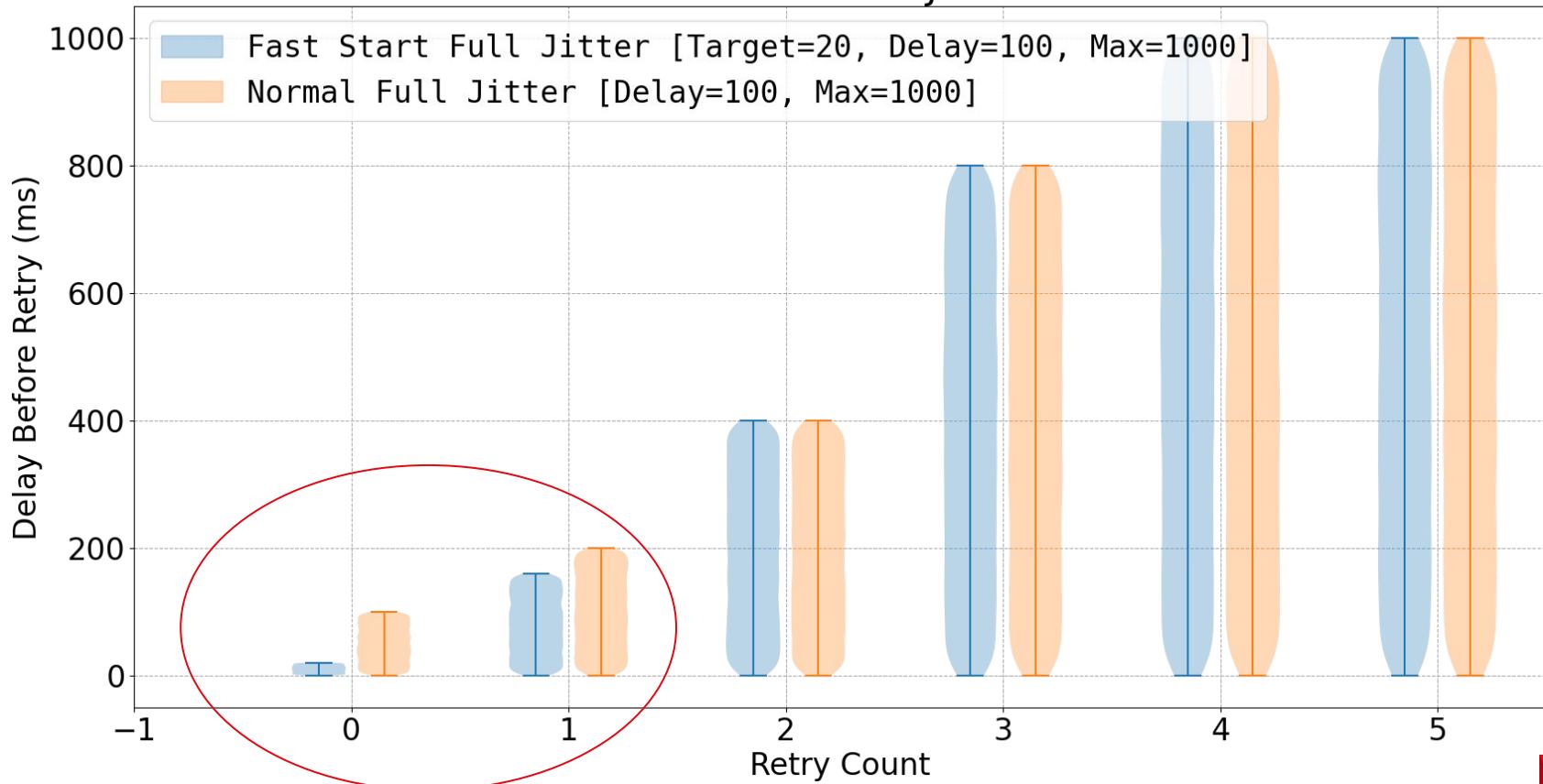
let  $R = \text{retry \#} \in [0, 1, 2, \dots \text{retry}_{\max} - 1]$

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

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

# Retries?

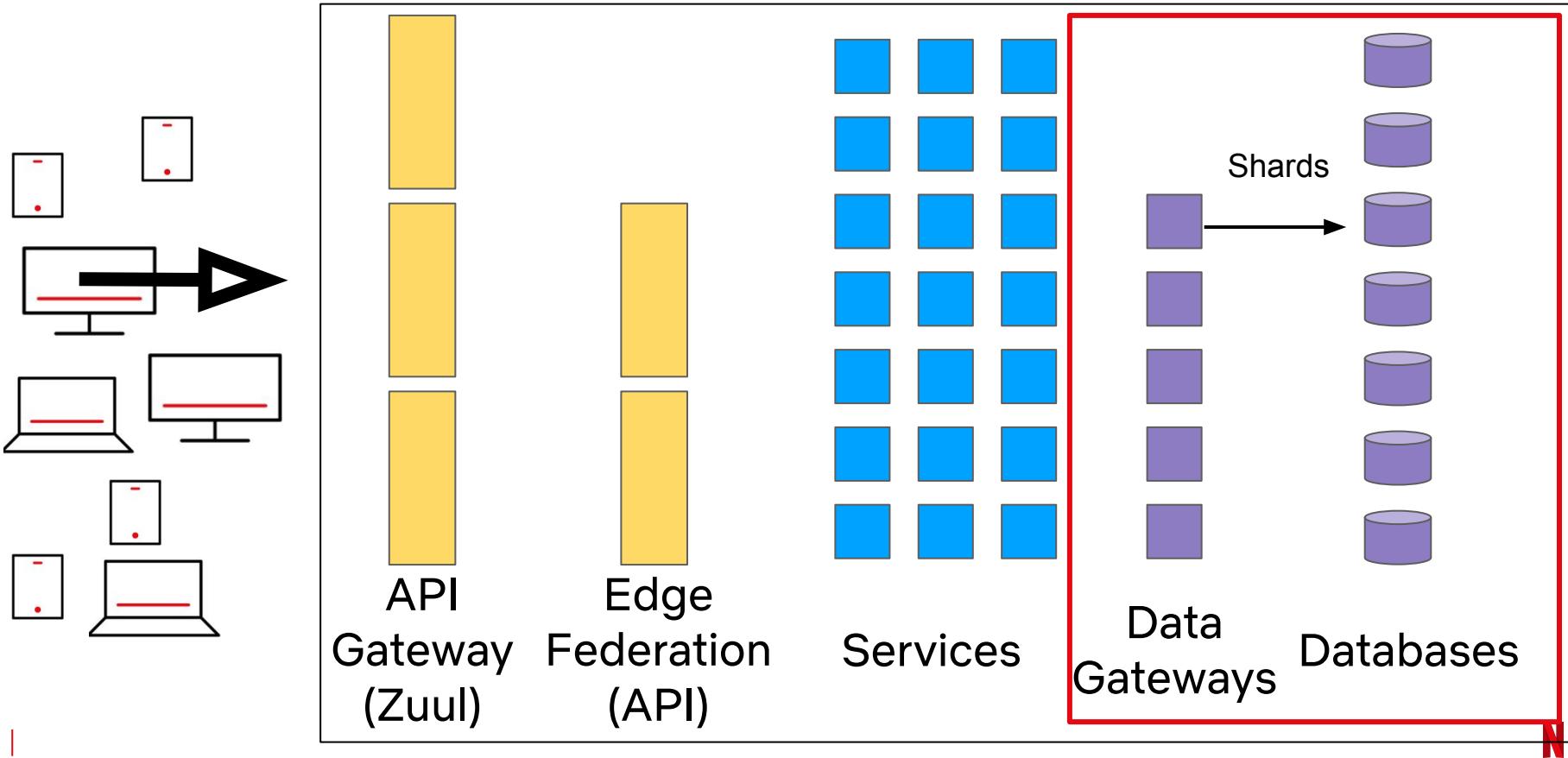
Fast Start Full Jitter



# **Stateful Resilience**

Capacity Planning  
Data Gateways  
Caching Wisely

# Microservices



## Capacity Plan Stateful

```
from service_capacity_modeling.capacity_planner import planner
from service_capacity_modeling.models.org import netflix

# Load up the Netflix capacity models
planner.register_group(netflix.models)

# Plan a cluster
plan = planner.plan(
    model_name="org.netflix.cassandra",
    region="us-east-1",
    desires=desires,
    simulations=1024,
    explain=True
)
```



## Strategy: Buffer, and lots of math

Least Regret Choice:

-----  
12 m5d.xlarge costing 8973.94

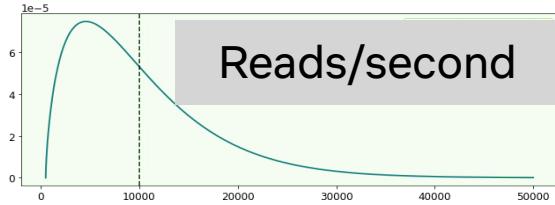
All Choices

-----  
{' 6 r5.xlarge': 4,  
 ' 6 r5d.large': 31,  
 ' 6 m5.2xlarge': 2,  
 ' 6 m5d.xlarge': 224,  
 ' 12 m5.xlarge': 132,  
 ' 12 m5d.xlarge': 277,  
 ' 24 m5.xlarge': 242,  
 ' 24 m5d.xlarge': 54,  
 ' 48 m5.xlarge': 55,  
 ' 48 m5d.xlarge': 2,  
 ' 96 m5.xlarge': 1}

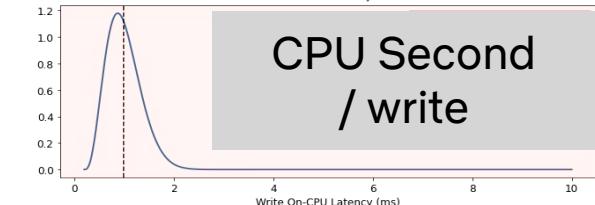
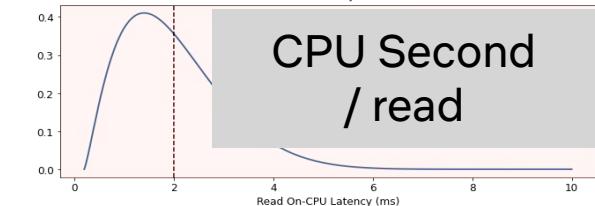
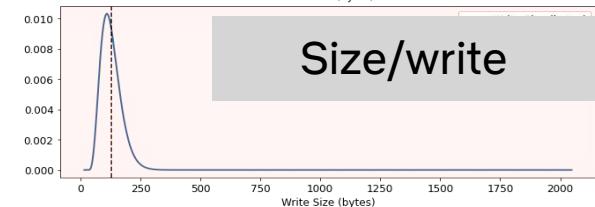
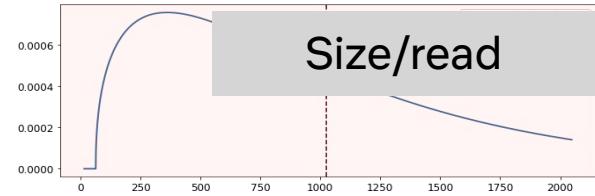
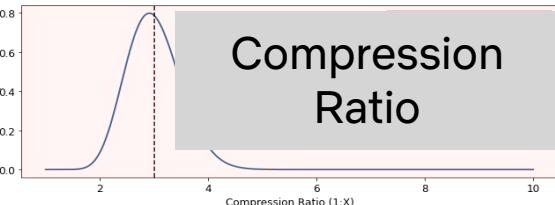
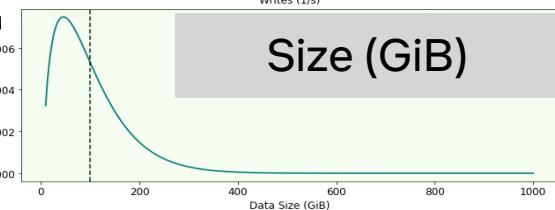
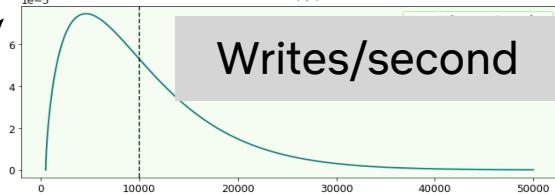
# Capacity Plan Stateful

# Optimally buy computers for each workload.

From the Human



From the Model



# Gateways Unlock Resilience

# chunked writes

```
x[1] ← d1{size = 1MiB, token = t1}  
x[2] ← d2{size = 1MiB, token = t1}  
x ← Commit{items = [d1, d2] } @ t1
```

# paginated reads

```
get(x) → {items = [d1], next = 2}  
get(x, prev = 2) → {items = [d2], next = ∅}  
x = d1 + d2
```



**Incremental**

# Gateways Unlock Resilience

## # chunked writes

```
x[1] ← d1{size = 1MiB, token = t1}  
x[2] ← d2{size = 1MiB, token = t1}  
x ← Commit{items = [d1, d2]} @ t1
```

# write  
 $x \leftarrow 1 @ t_1$   
 $\text{get}(x) \rightarrow 1$

# update  
 $x \leftarrow 2 @ t_2$   
# hedge  
 $x \leftarrow 1 @ t_1$

## # paginated reads

```
get(x) → {items = [d1], next = 2}  
get(x, prev = 2) → {items = [d2], next = ∅} # update visible  
x = d1 + d2  
get(x) → 2
```



**Incremental**

**Idempotent**

# Gateways Unlock Resilience

# chunked writes

```
x[1] ← d1{size = 1MiB, token = t1}  
x[2] ← d2{size = 1MiB, token = t1}  
x ← Commit{items = [d1, d2]} @ t1
```

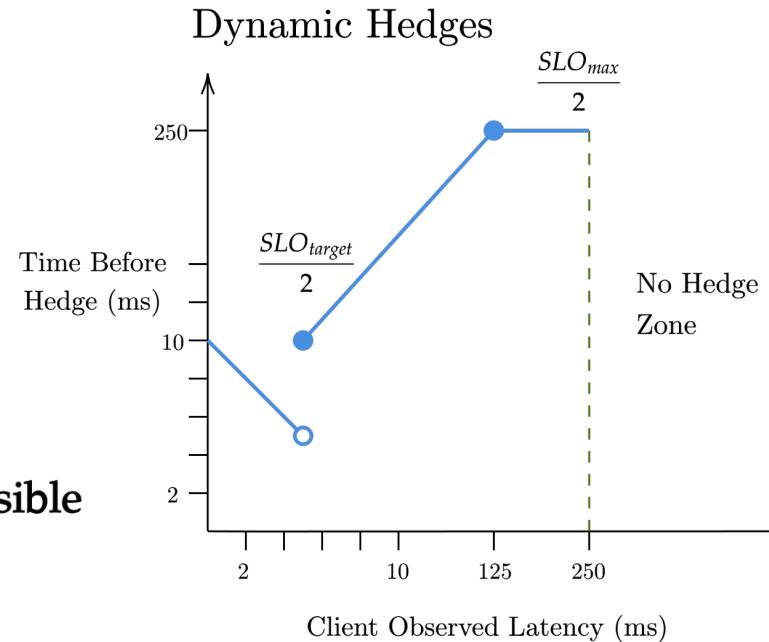
# paginated reads

```
get(x) → {items = [d1], next = 2}  
get(x, prev = 2) → {items = [d2], next = ∅}  
x = d1 + d2
```

# write  
 $x \leftarrow 1 @ t_1$   
 $\text{get}(x) \rightarrow 1$

# update  
 $x \leftarrow 2 @ t_2$   
# hedge  
 $x \leftarrow 1 @ t_1$

# update visible  
 $\text{get}(x) \rightarrow 2$



Incremental

Idempotent

Retriable

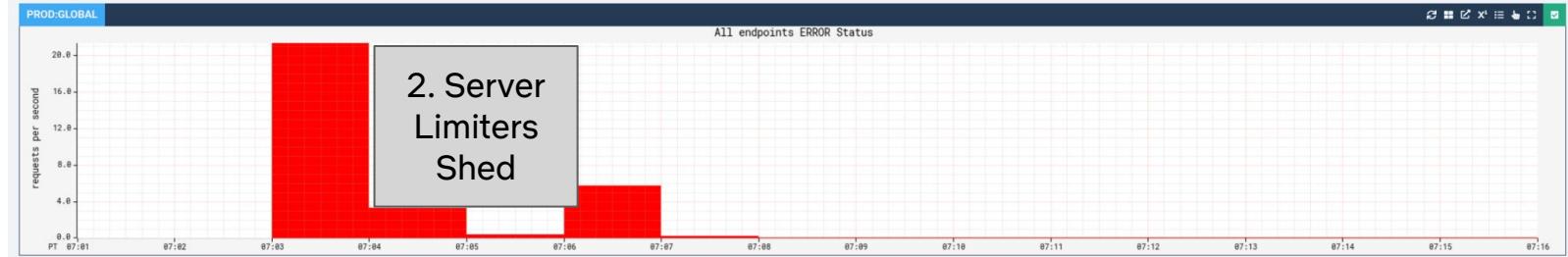
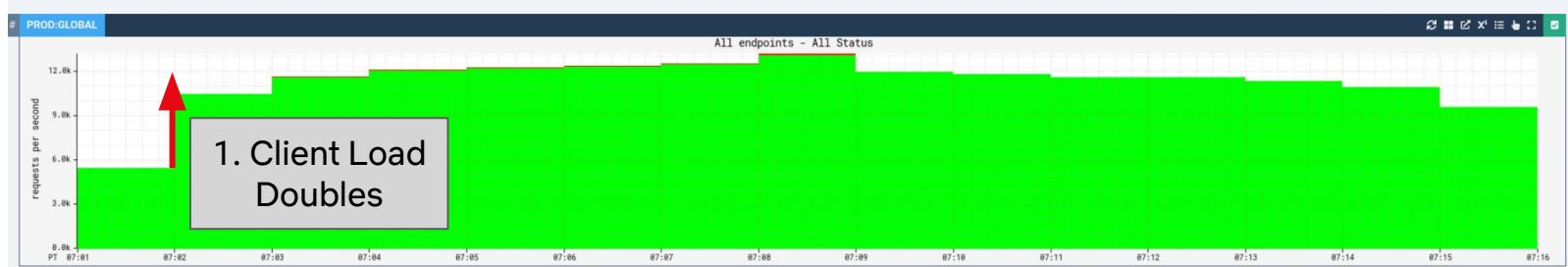
# Load Shedding and Retries

# Prefer limiting on server, clients implement basic safety



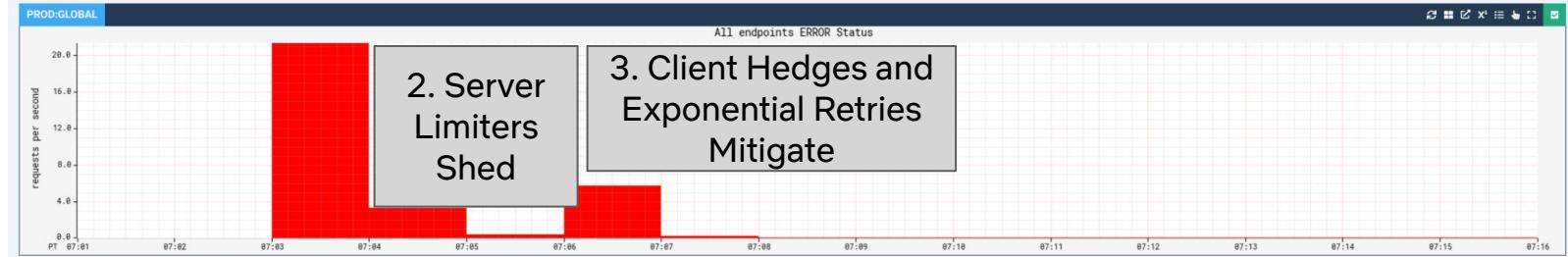
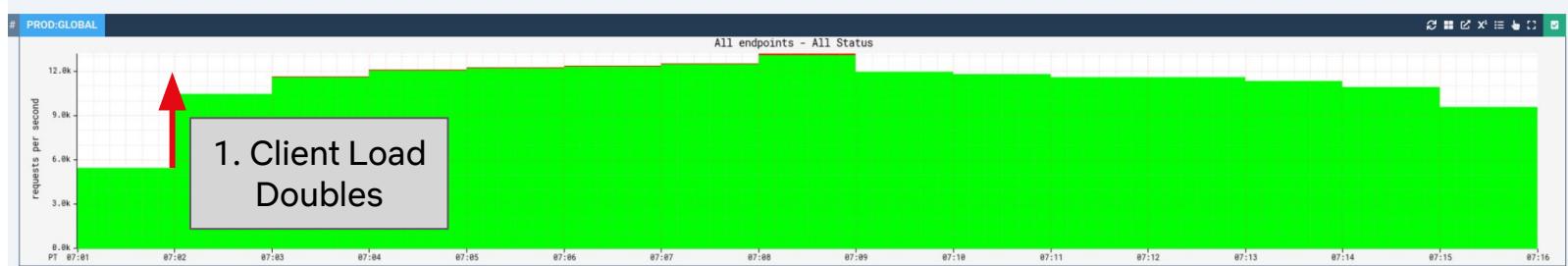
# Load Shedding and Retries

# Prefer limiting on server, clients implement basic safety



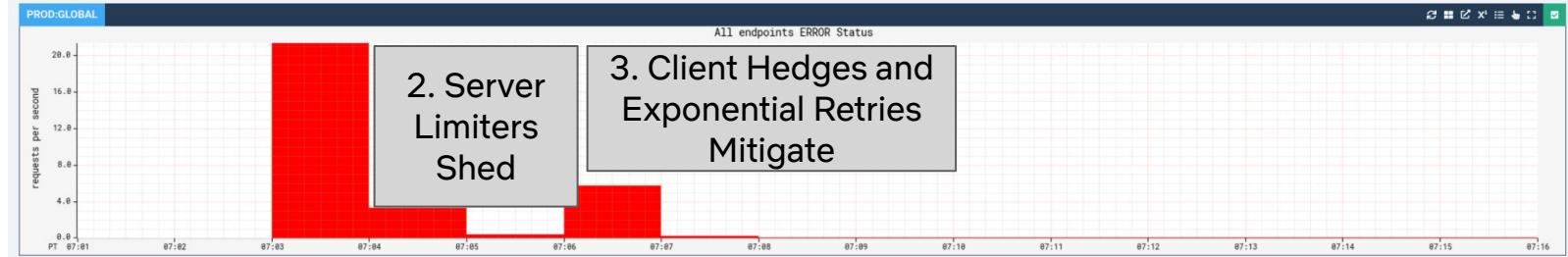
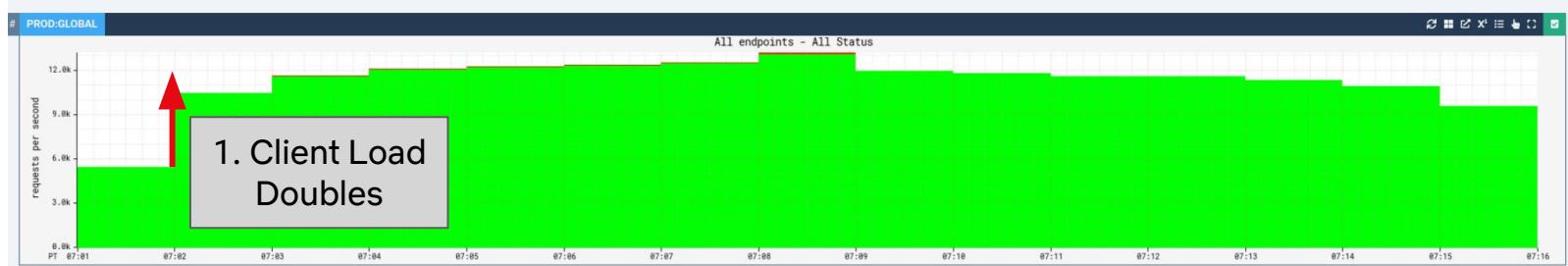
# Load Shedding and Retries

# Prefer limiting on server, clients implement basic safety



# Load Shedding and Retries

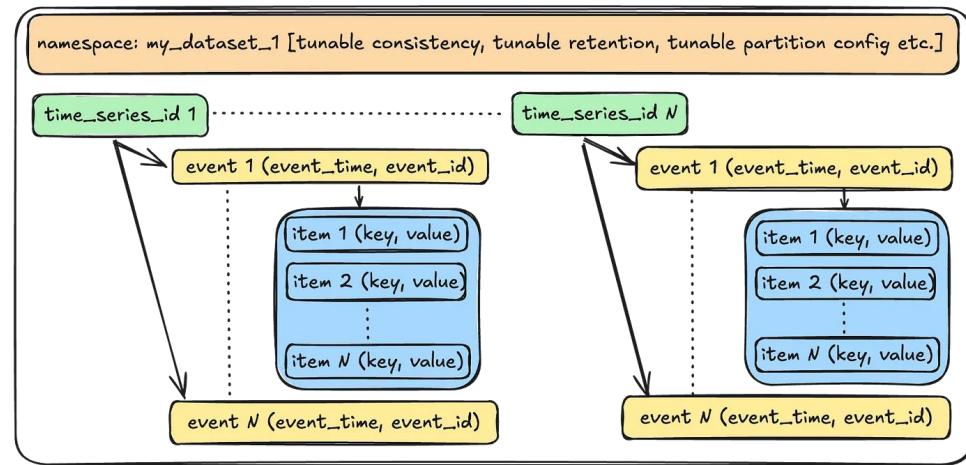
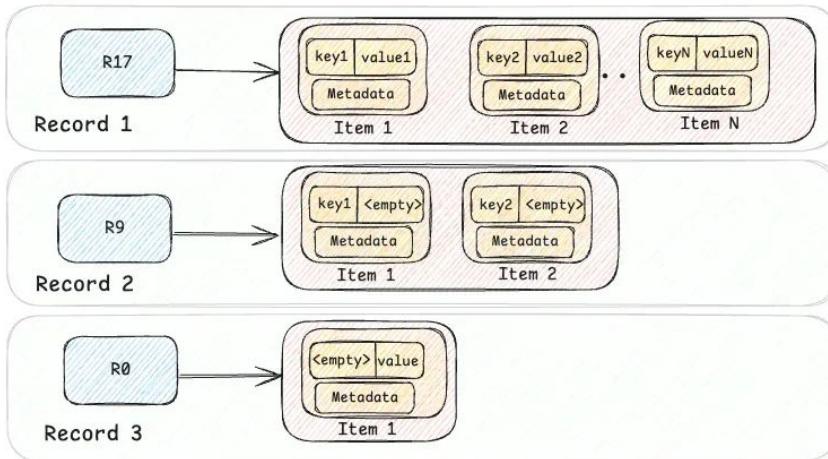
# Prefer limiting on server, clients implement basic safety



# Key-Value



Record ID (ID) → Item[Key -> Value]...

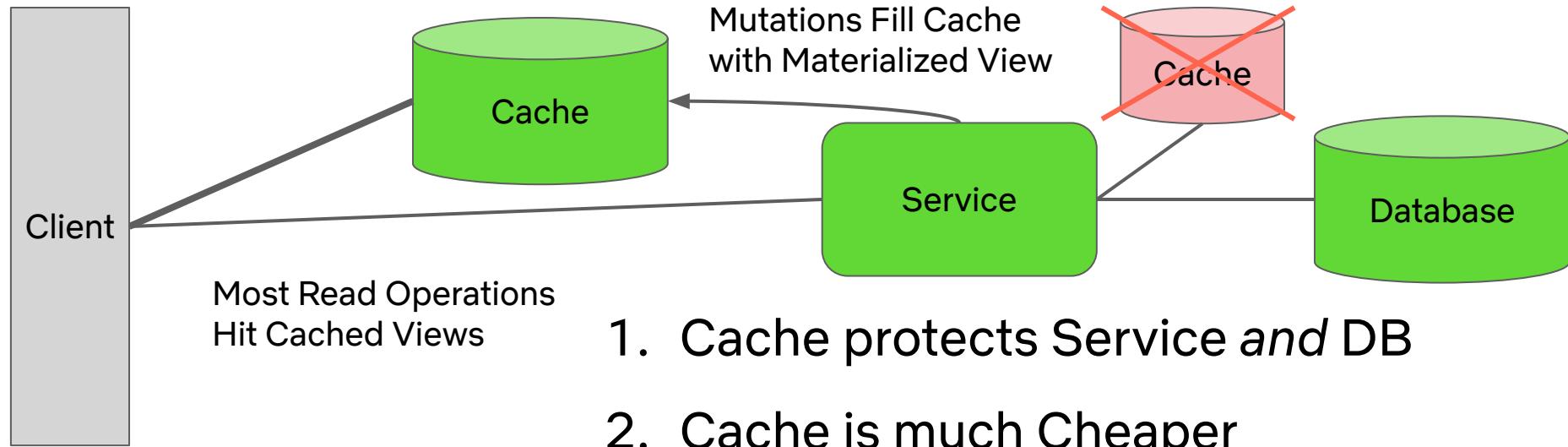


# Time-Series



## Cache In Front of Services

# Cache your service, not your database!

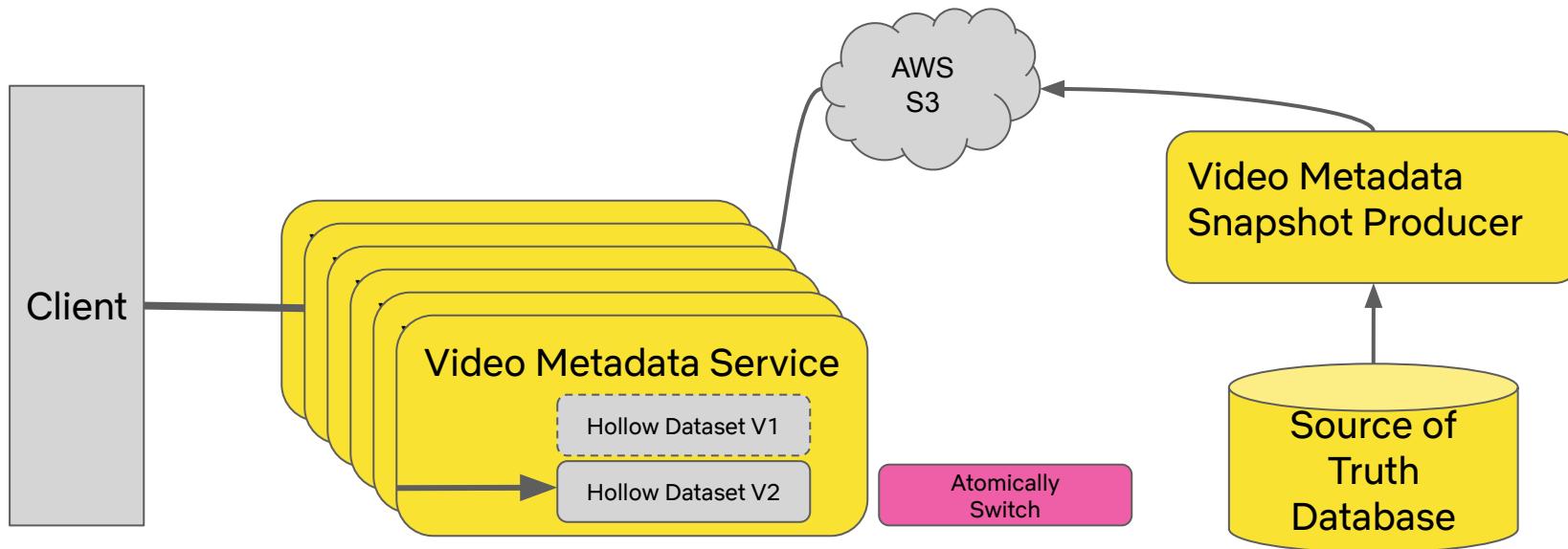


1. Cache protects Service and DB
2. Cache is much Cheaper
3. Cache can be Replicated



## Total Cache What You Can

Replicated total near caches like Netflix  
Hollow handle Load increases easily!



# Test Constantly

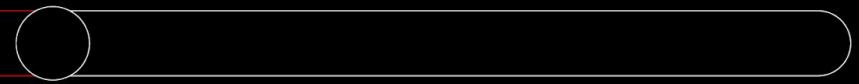
Properties Tests

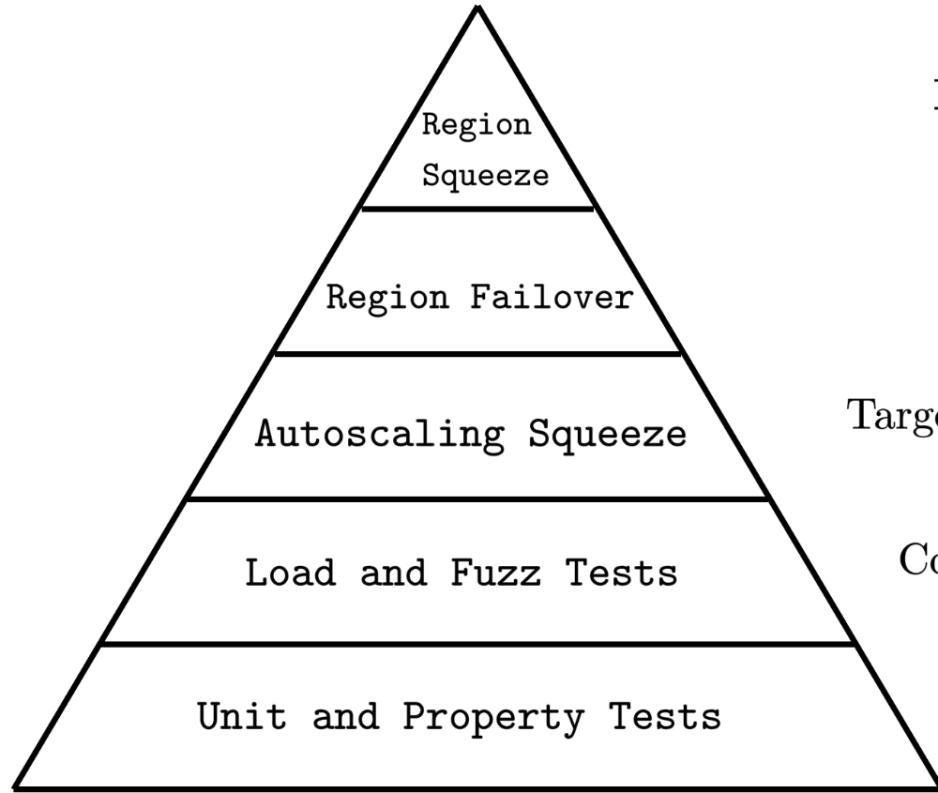


Load Tests



Failover Tests





**De-risk the Unknown** with Full Global Traffic

Full System **Integration Tests**, Verify Buffers

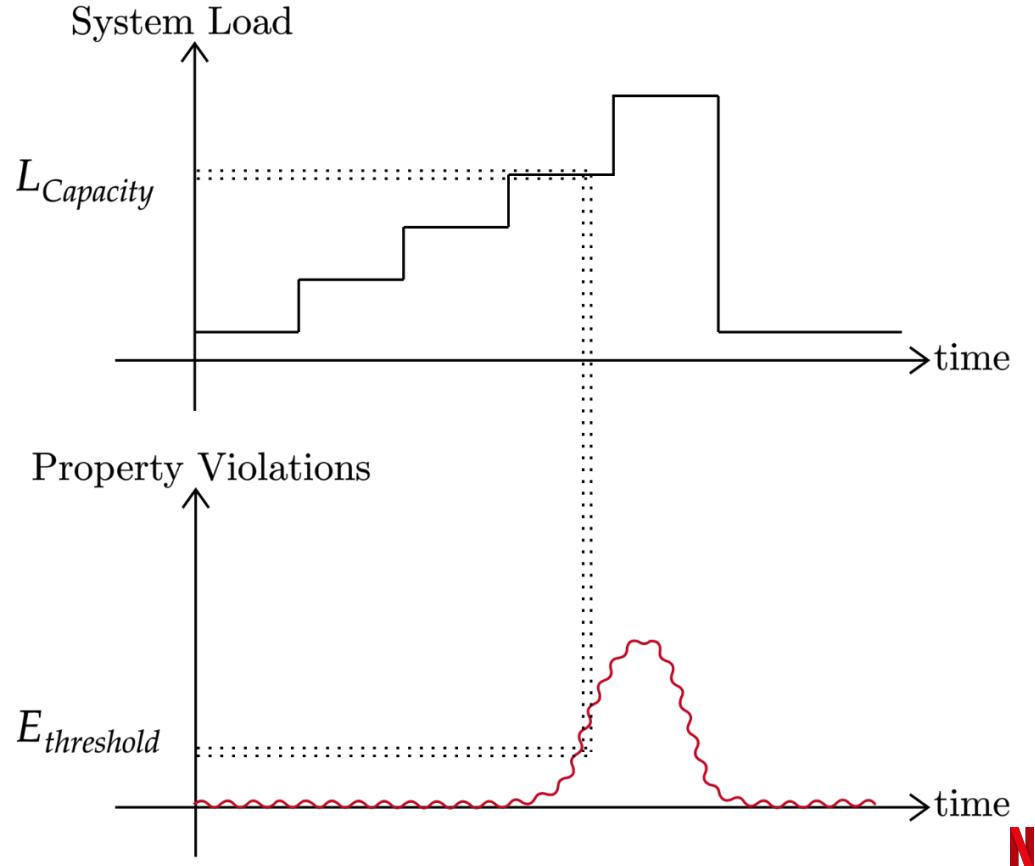
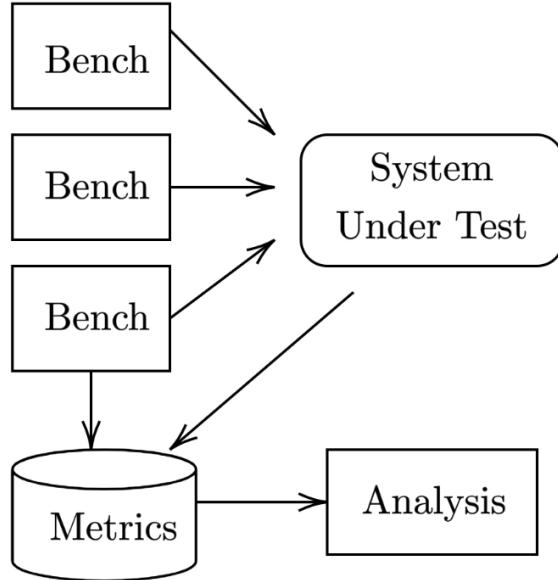
Targeted **Scaling Tests** with **Property Assertions**

Controlled **Benchmarks** and **Randomized Fuzz**

Reliable **Unit** and **Randomized Properties**

# Example Load Test

```
{  
    write_rate = [0, 10k],  
    write_bytes = [1, 4KiB], ...  
}
```



## Manage Traffic Demand

Full-Active Control of Demand

- Global **traffic shaping**
- Prioritization at the Edge
- **Fallbacks!**

Understand Traffic Flows

- Uneven Service Flows
- Different **Traffic Priorities**

Manage Quality of Service

- **Slow** often **better than Down**

## Balance Compute Supply

Capacity Planning

- Compute and Workload Analysis
- Intentional **Buffers**
- **Pre-scale** when you can

Get Out of Trouble

- **Autoscale**
- **Shed Load**
- Stateful can be reliable too!

Manage Quality of Service

- **Prioritized** Shedding
- **Stale** often **better than "Down"**

# Thank You.

Joseph Lynch

[josephl@netflix.com](mailto:josephl@netflix.com)



[jolynch.github.io](https://jolynch.github.io)

