

How Netflix Shapes our Fleet for Efficiency and Reliability

Argha C and Joey Lynch

N

Speaker

Argha C



Staff Software Engineer
Infrastructure at Netflix

Focus on Compute and Network
problems.

<https://www.linkedin.com/in/argha-c/>



Speaker

Joseph Lynch



Principal Software Engineer
Platform Engineering at Netflix

Database shepherd, compute optimizer,
distributed system nerd

<https://jolynch.github.io/>



Agenda

Goal: Efficiency *and* Reliability

Compute **Supply**

— Capacity Planning → Fleet Planning

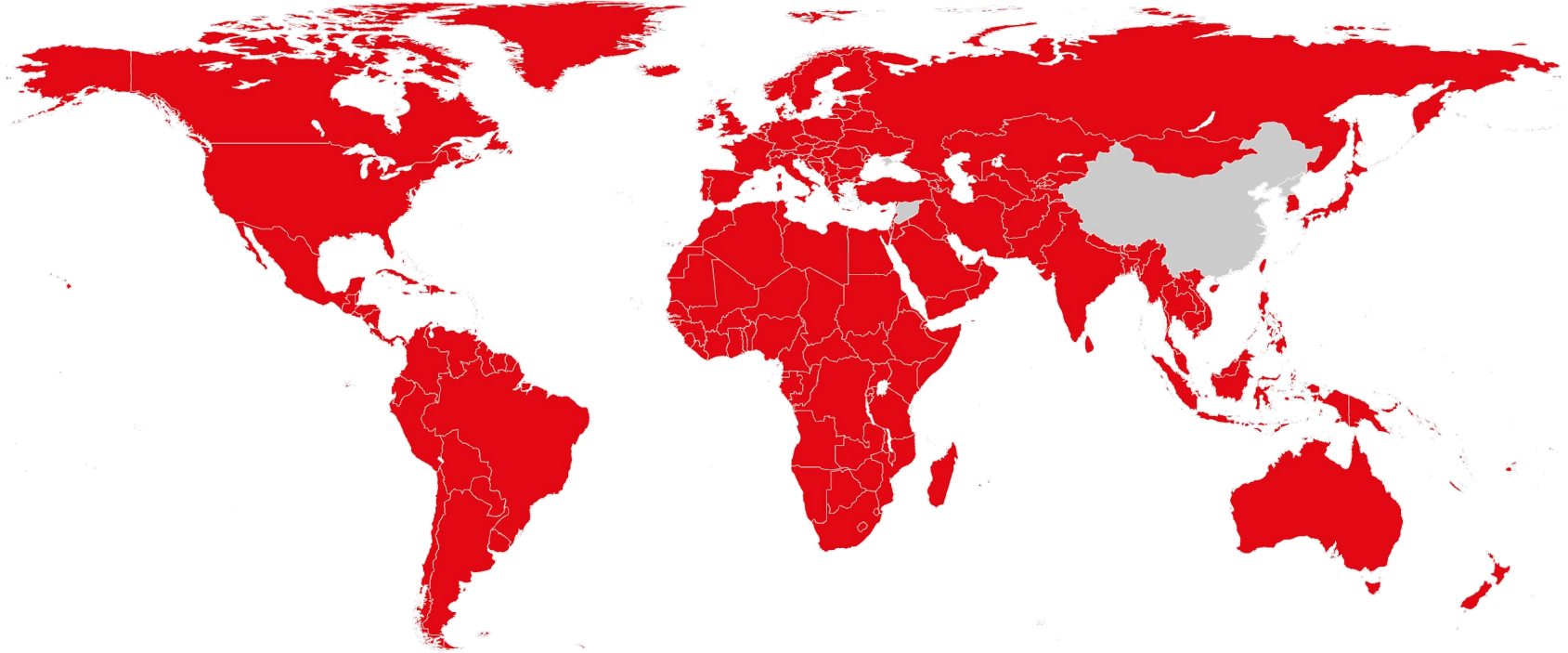
Workload **Demand**

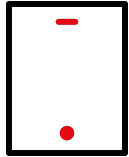
— Workload Shape → Traffic Shaping

Balance Supply with Demand

Reliability Techniques **Allow Risk**

Problem - Global





Mobile

Diverse Device Capabilities

Weaker Network

Android and IOs

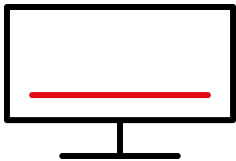


PC

Diverse Device Capabilities

More Stable Network

Medium Resolution



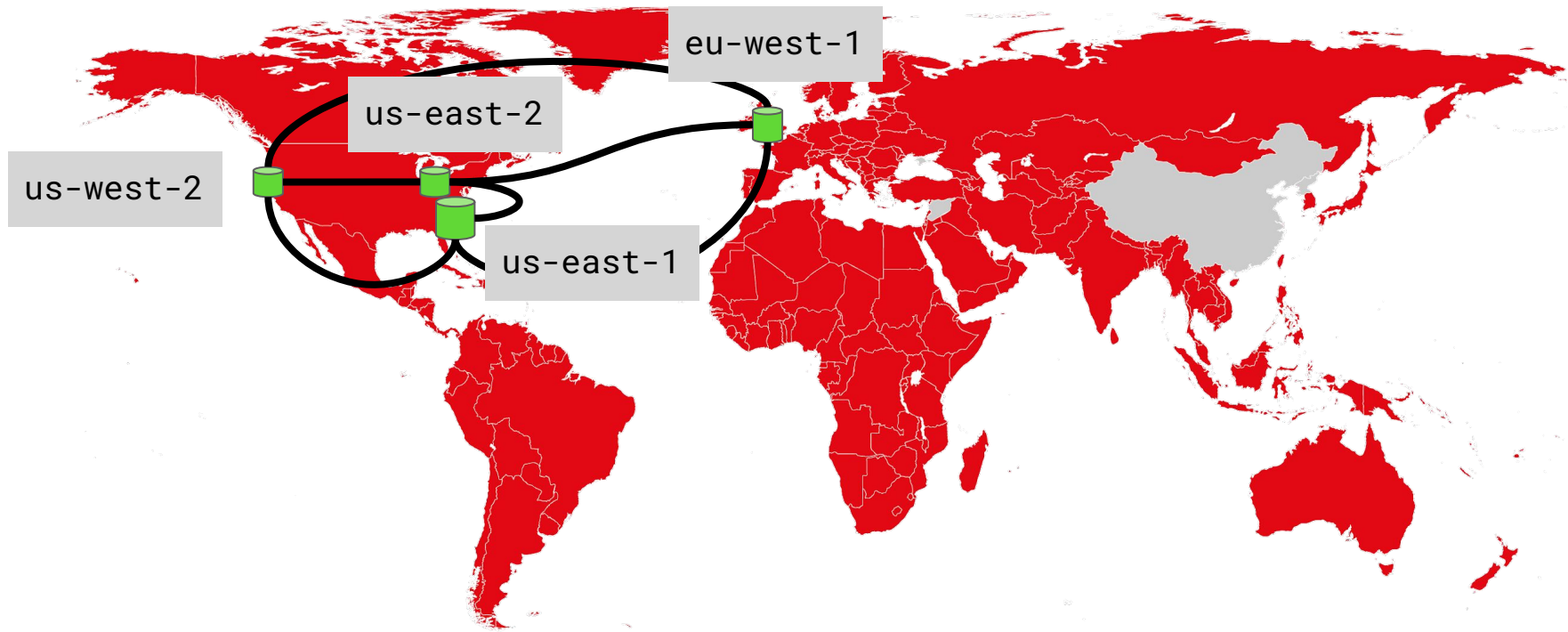
TV

Diverse Device Capabilities

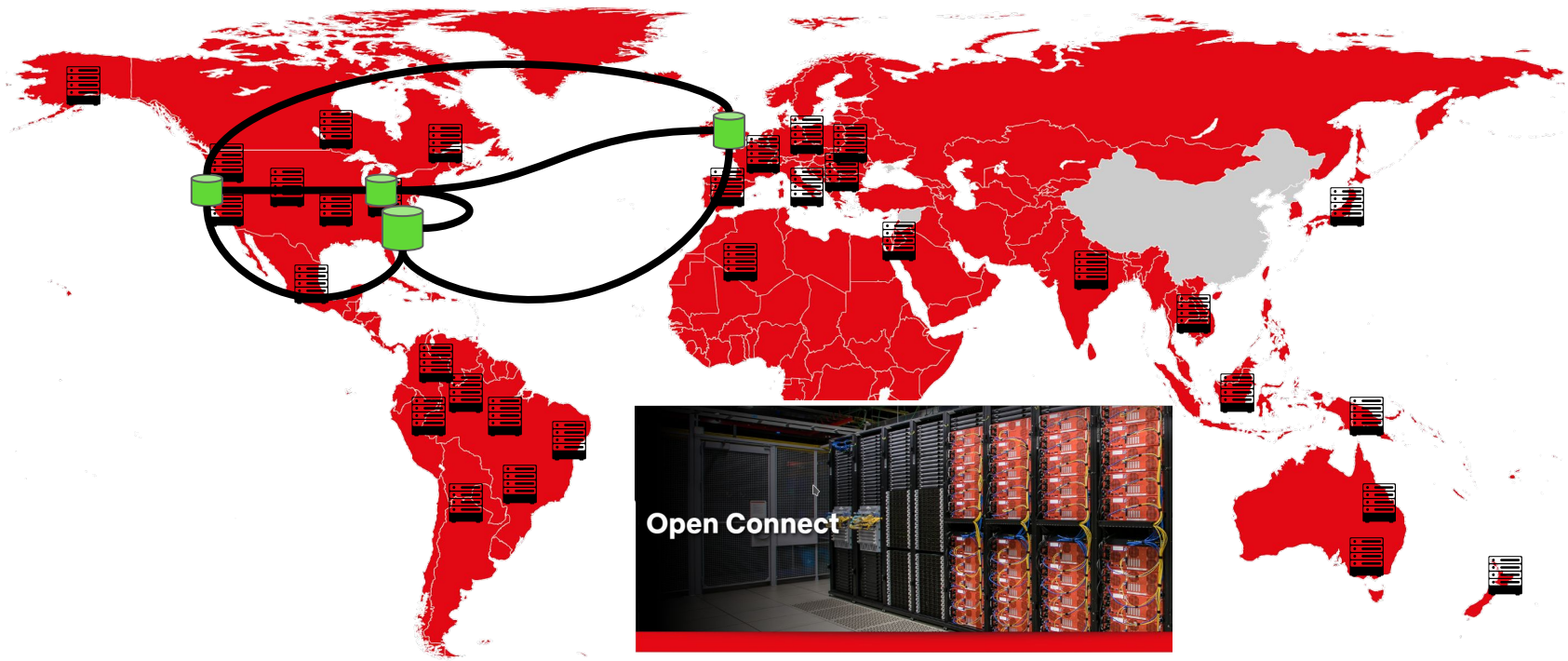
More Stable Network

High Resolution

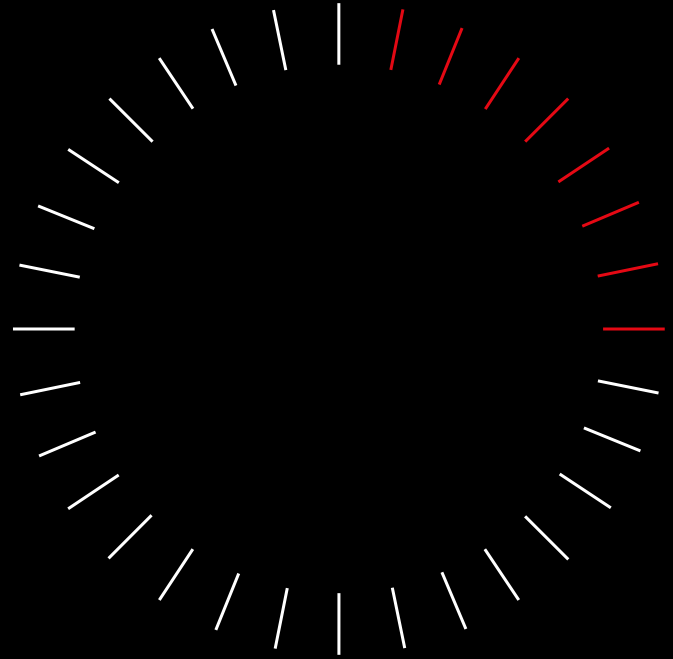
Solution - Global Control Plane



Solution - Global Data Plane



Efficiency *and* Reliability



Efficiency

Workloads create **Business Value**

Efficiency

Workloads create **Business Value**

Workloads have **Cost**

Efficiency

Workloads create **Business Value**

Workloads have **Cost**

Workload **Cost** \propto **Resource Usage**

Efficiency

Workloads create **Business Value**

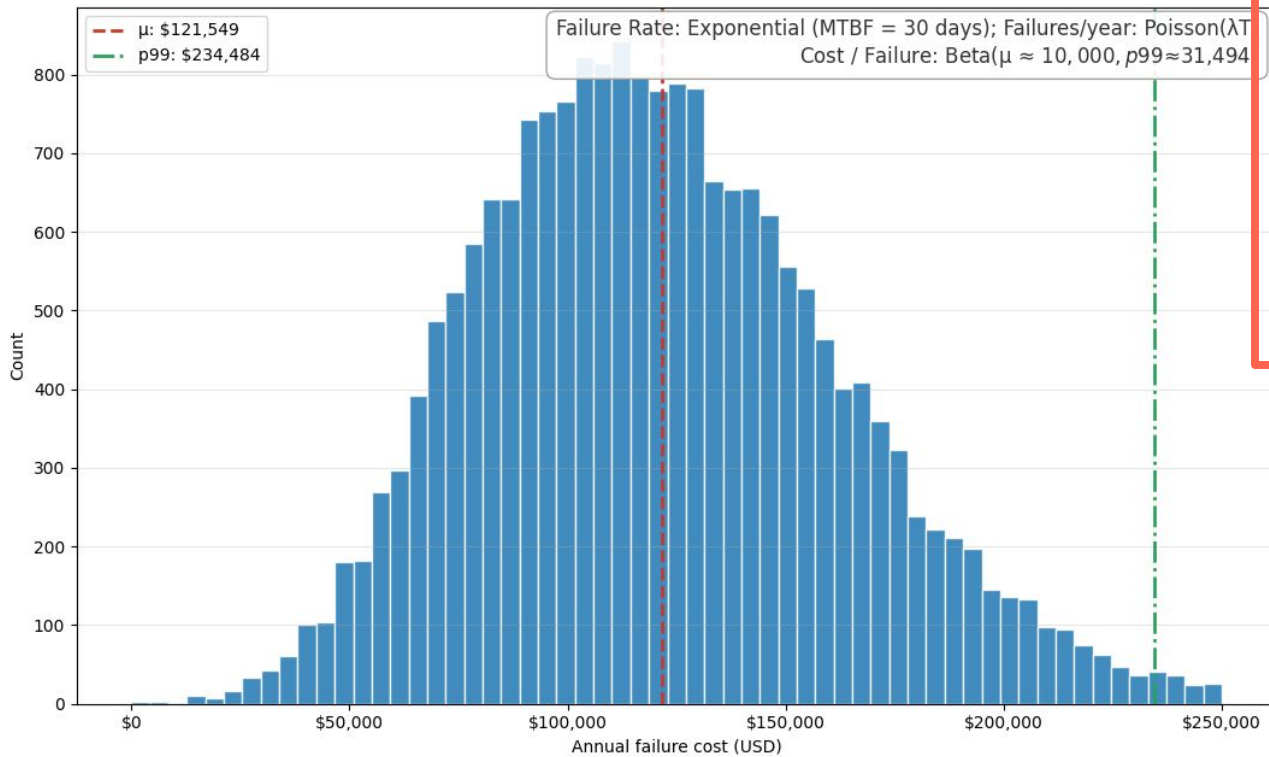
Workloads have **Cost**

Workload **Cost** \propto **Resource Usage**

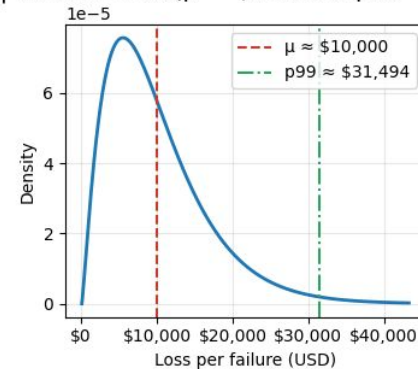
Workload **Failures** have **Cost**

Failures have cost!

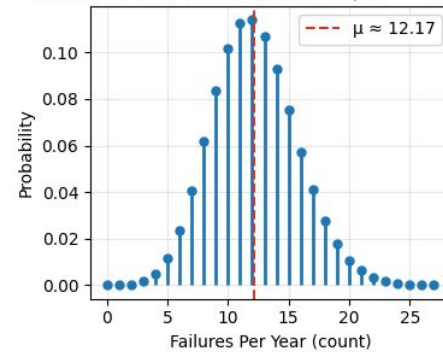
Annual Failure Cost Distribution



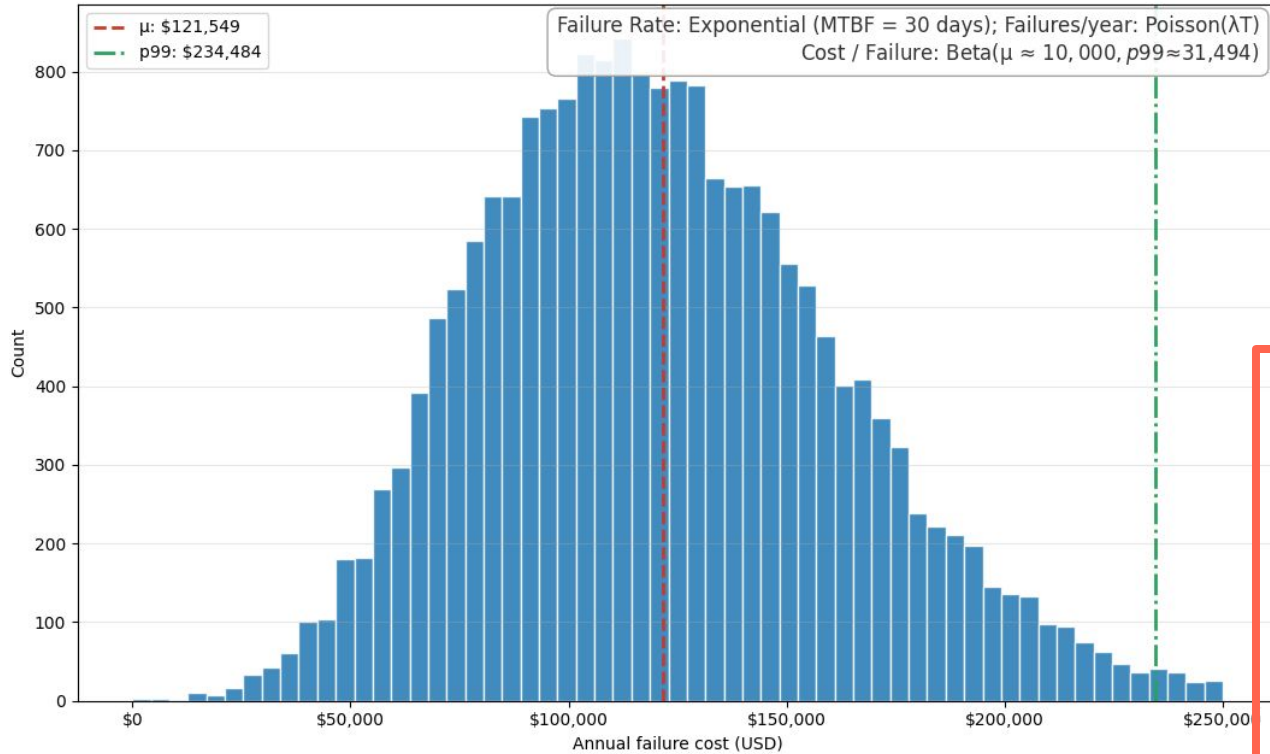
Loss | Failure: Beta($\mu \approx \$10,000$, $p_{99} \approx \$31,494$)



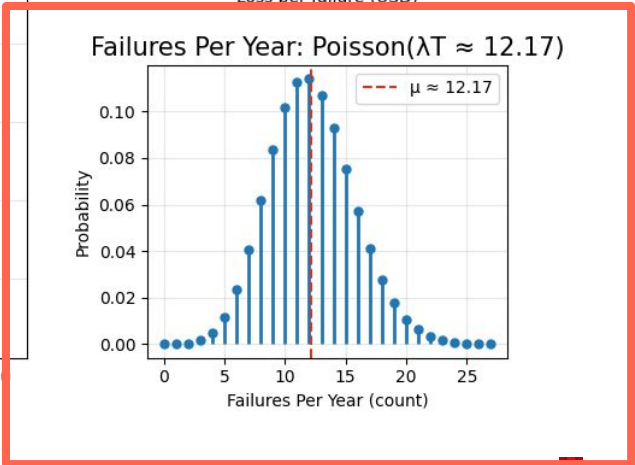
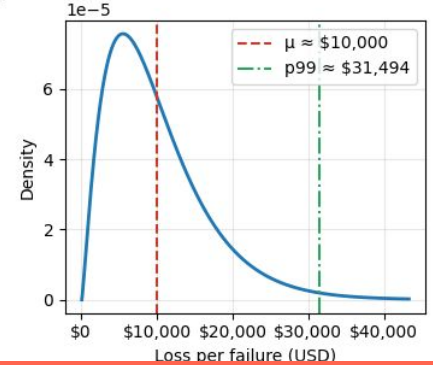
Failures Per Year: Poisson($\lambda T \approx 12.17$)



Annual Failure Cost Distribution

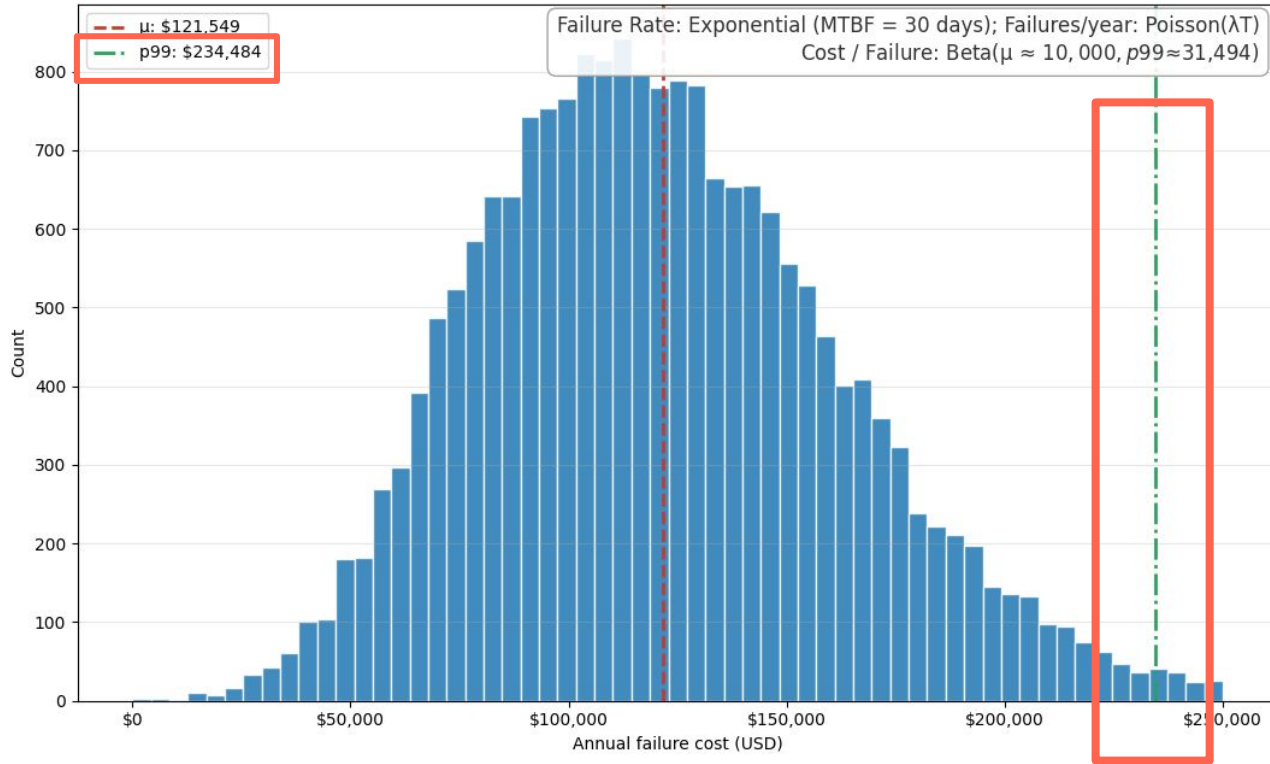


Loss | Failure: Beta($\mu \approx \$10,000, p99 \approx \$31,494$)

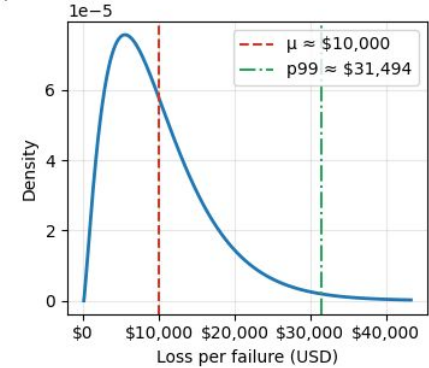


Failures happen at some rate!

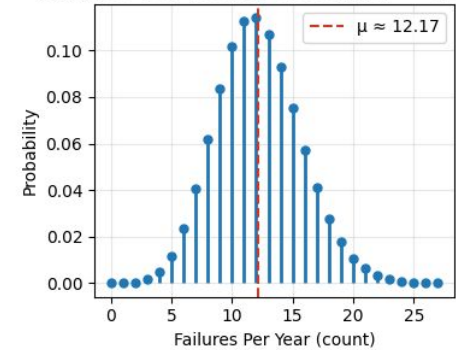
Annual Failure Cost Distribution



Loss | Failure: Beta($\mu \approx \$10,000$, $p99 \approx \$31,494$)



Failures Per Year: Poisson($\lambda \approx 12.17$)

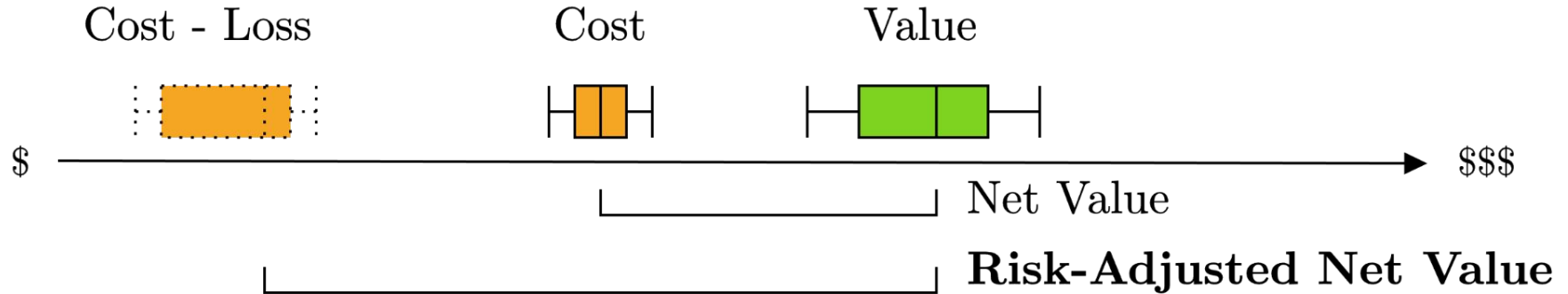


You probably care about the tail



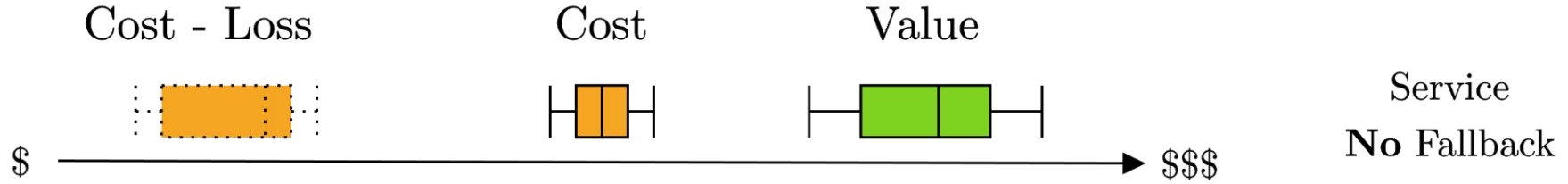
What is Efficiency?

Balance Value, Cost, and Risk



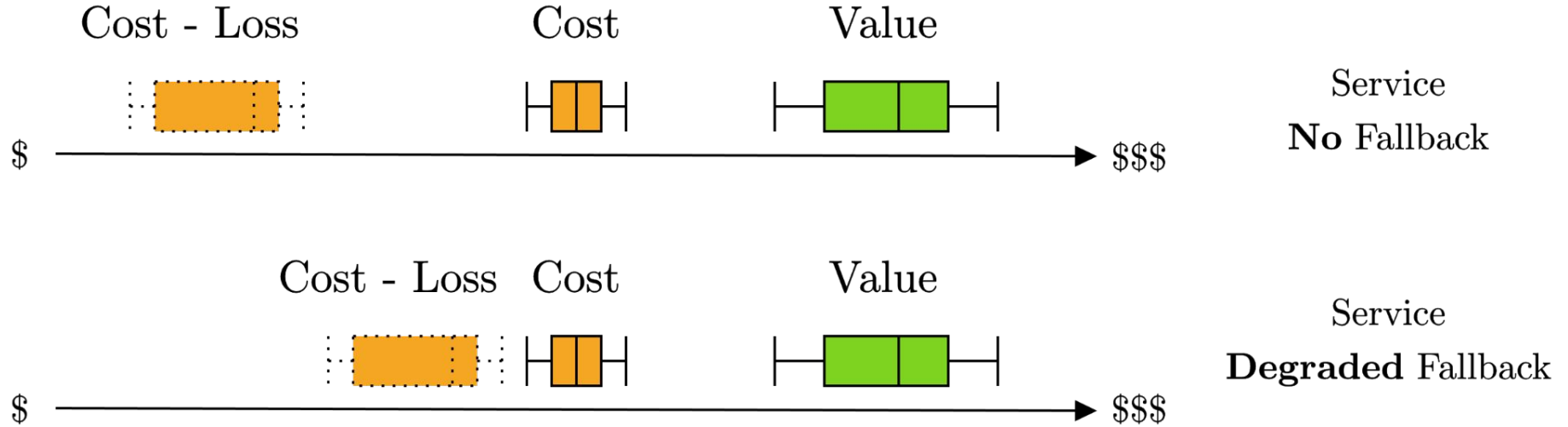
What is Efficiency?

Balance Value, Cost, and Risk



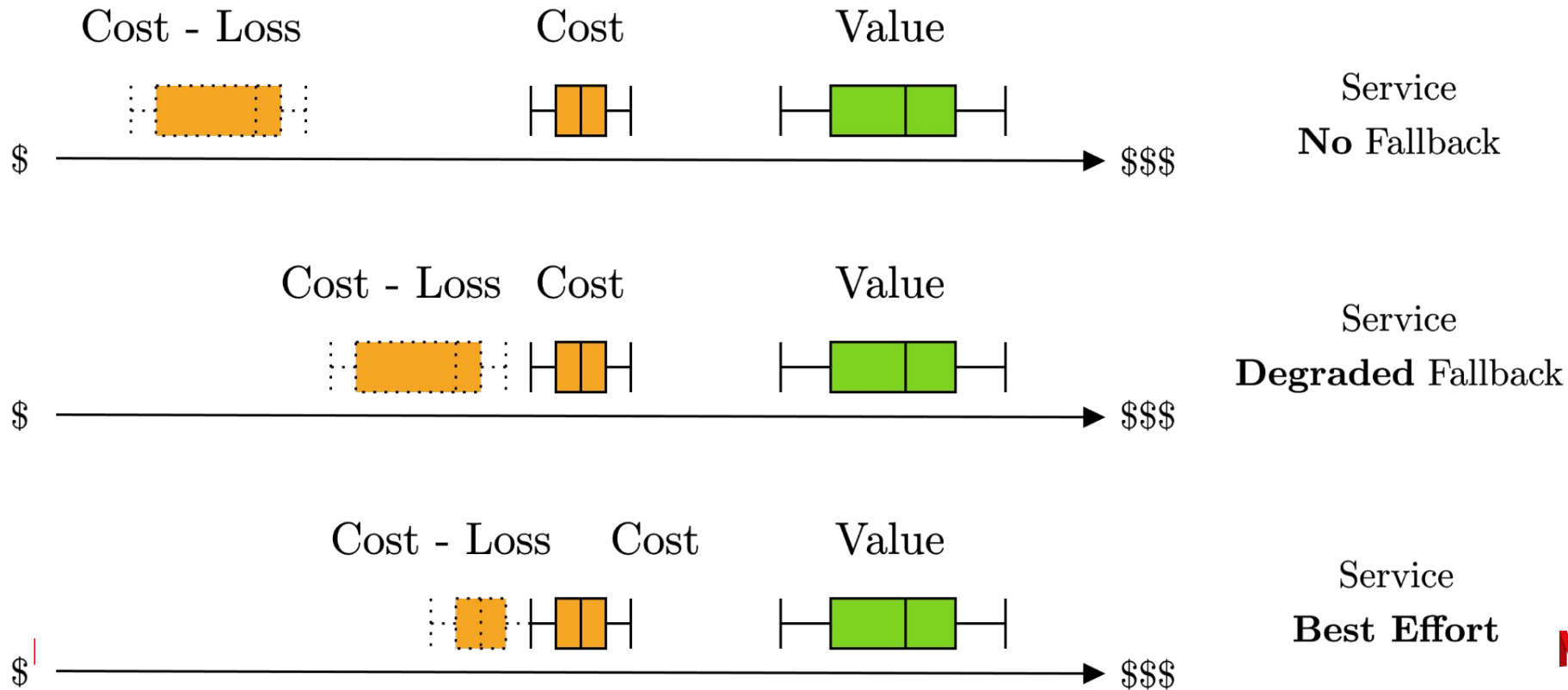
What is Efficiency?

Balance Value, Cost, and Risk



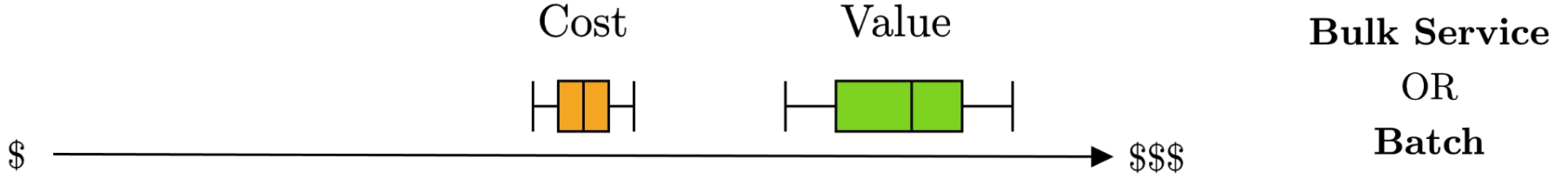
What is Efficiency?

Balance Value, Cost, and Risk



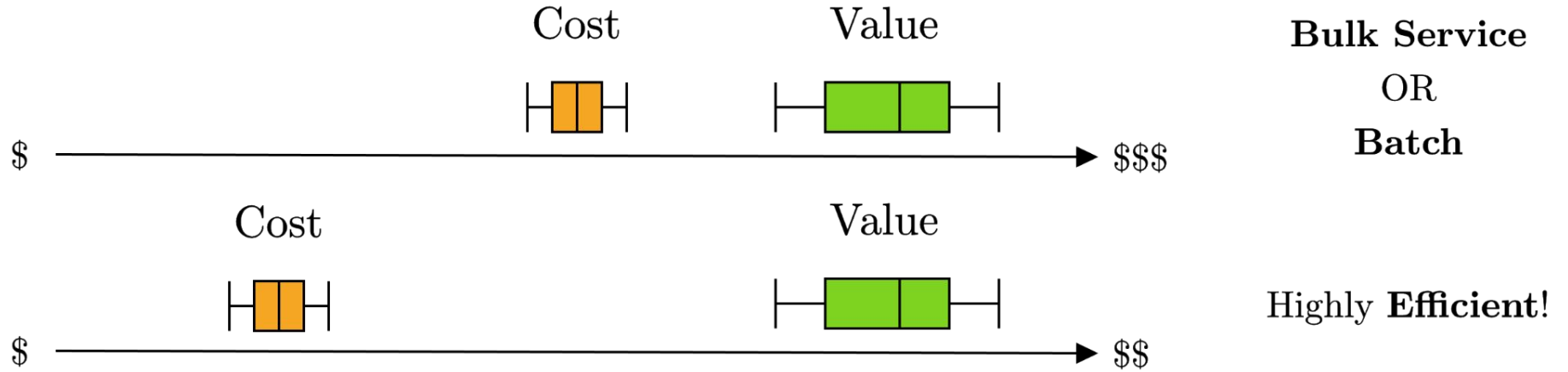
What is Efficiency?

Balance Value, Cost, and Risk



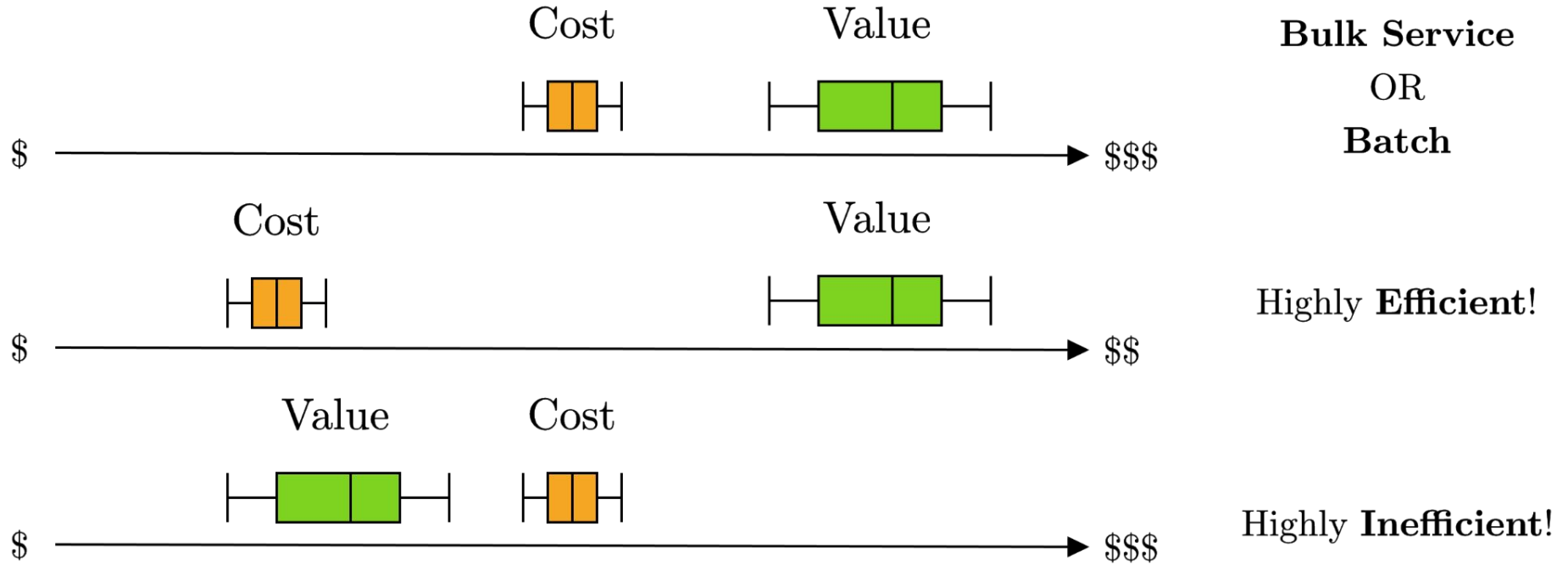
What is Efficiency?

Balance Value, Cost, and Risk



What is Efficiency?

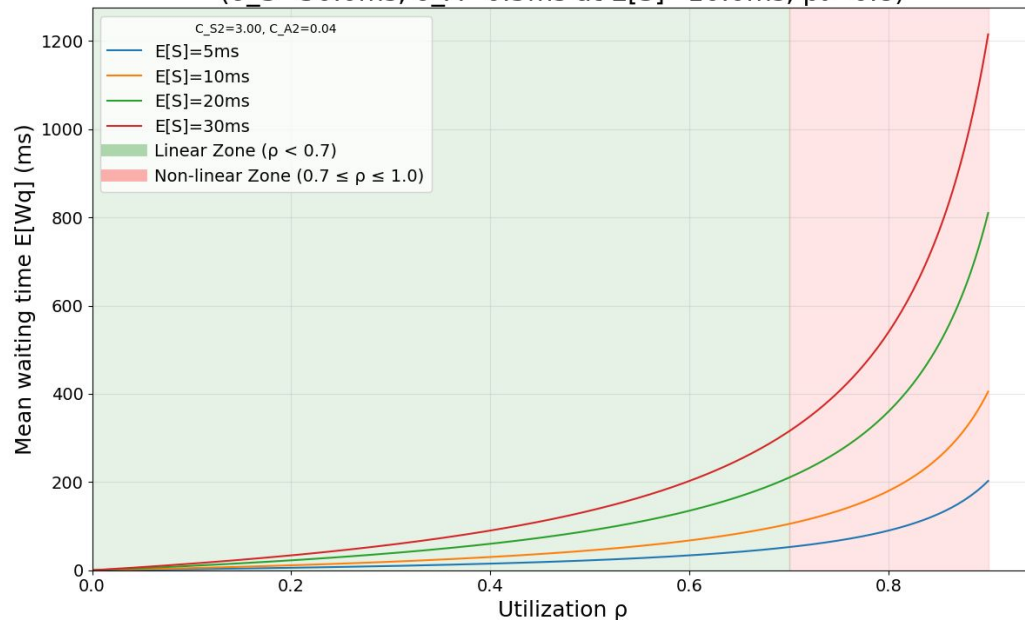
Balance Value, Cost, and Risk



System Efficiency

$$\mathbb{E}[Delay] \approx \frac{\rho}{1 - \rho} \times \mathbb{E}[S] \times \left(\frac{C_S^2 + C_A^2}{2} \right)$$

Kingman's Approximation of Queueing Delay
($\sigma_S=30.0\text{ms}$, $\sigma_A=0.5\text{ms}$ at $\mathbb{E}[S]=10.0\text{ms}$, $\rho_0=0.8$)



System Efficiency

Raising Utilization is just one knob

$$\mathbb{E}[Delay] \approx \frac{\rho}{1 - \rho} \times \mathbb{E}[S] \times \left(\frac{C_S^2 + C_A^2}{2} \right)$$

Kingsman's Approximation

$\rho =$ system load (utilization)

$\mathbb{E}[S] =$ avg svc time

$C_S^2 =$ variation of svc time

$C_A^2 =$ variation of arrival rate

Make the Service Faster, a.k.a

Performance Engineering!



System Efficiency

Raising Utilization is just one knob

$$\mathbb{E}[\text{Delay}] \approx \frac{\rho}{1 - \rho} \times \mathbb{E}[S] \times \left(\frac{C_S^2 + C_A^2}{2} \right)$$

Kingsman's Approximation

$$\begin{aligned} \rho &= \text{system load (utilization)} \\ \mathbb{E}[S] &= \text{avg svc time} \\ C_S^2 &= \text{variation of svc time} \\ C_A^2 &= \text{variation of arrival rate} \end{aligned}$$

Make the Service *Consistently* Fast, a.k.a

Performance Engineering!



System Efficiency

Raising Utilization is just one knob

$$\mathbb{E}[Delay] \approx \frac{\rho}{1 - \rho} \times \mathbb{E}[S] \times \left(\frac{C_S^2 + C_A^2}{2} \right)$$

Kingsman's Approximation

$$\begin{aligned} \rho &= \text{system load (utilization)} \\ \mathbb{E}[S] &= \text{avg svc time} \\ C_S^2 &= \text{variation of svc time} \\ C_A^2 &= \text{variation of arrival rate} \end{aligned}$$

Balance Load (Choice-of-N)

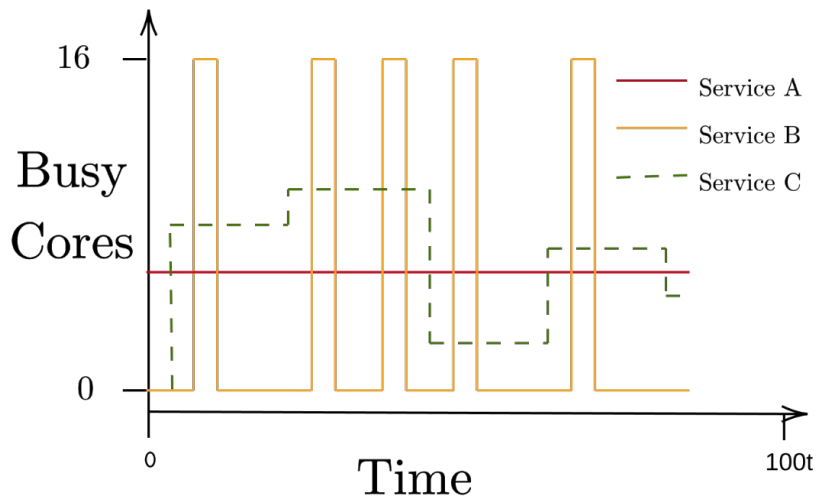
Spread Load (Token Buckets)



System Efficiency

Utilization is often a Lie

30% Utilization of a 16-Core System



Kingsman's Approximation

$\rho =$ system load (utilization)

$\mathbb{E}[S] =$ avg svc time

$C_S^2 =$ variation of svc time

$C_A^2 =$ variation of arrival rate

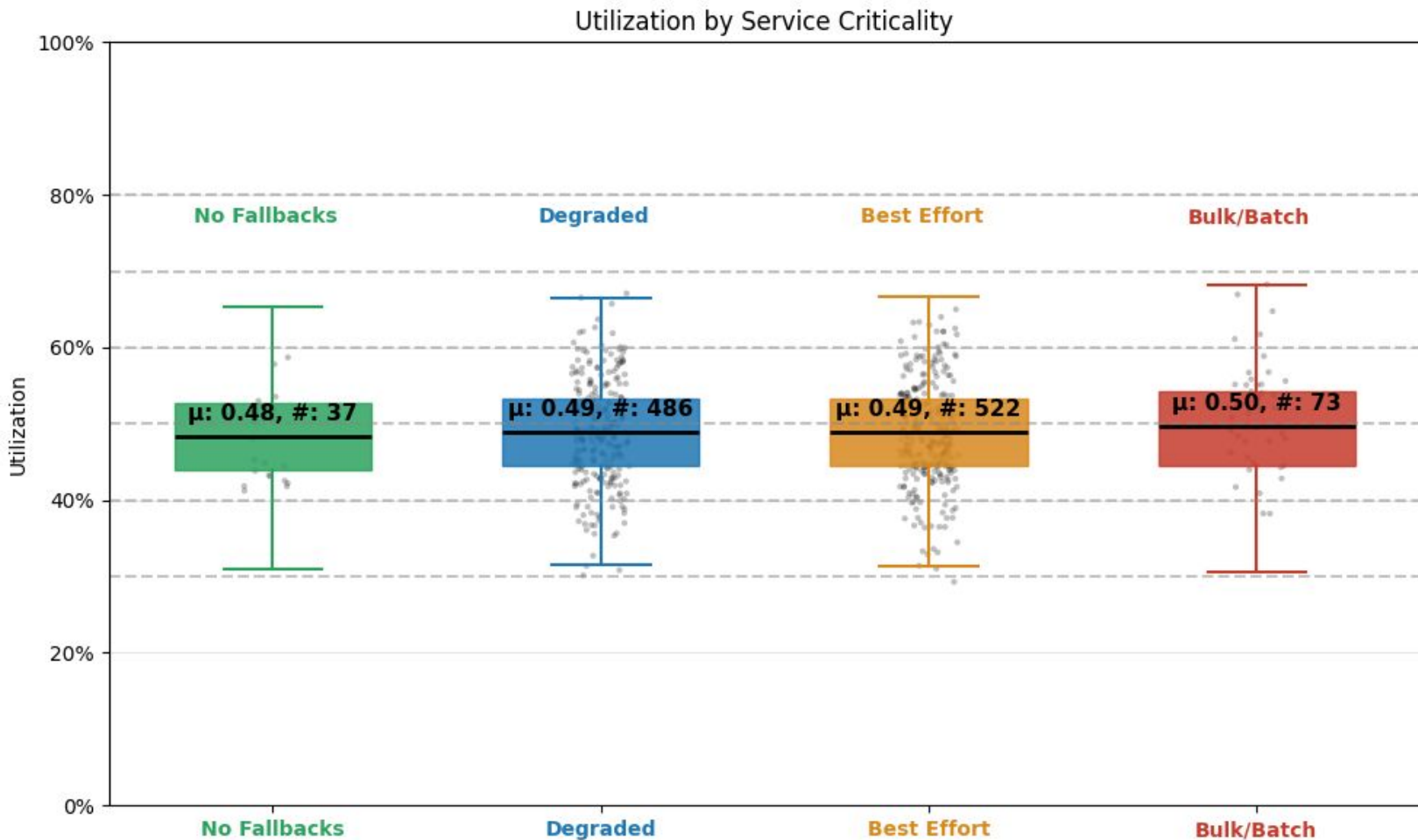
Instead of Raising Utilization, ask

Where do we spend our dollars?

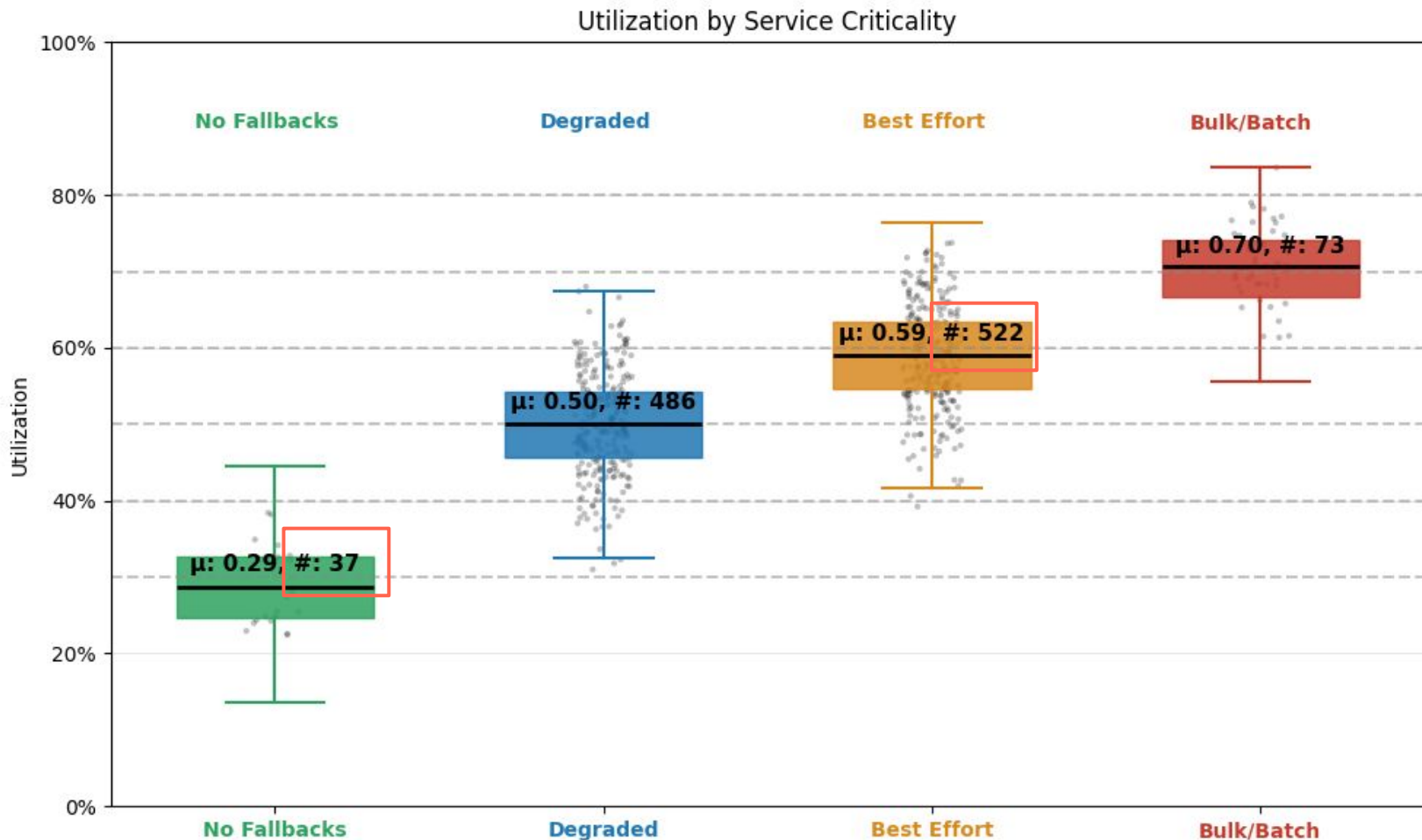
Do those dollars mitigate risk?

What lever do I pull to reduce cost?

Not Efficient Fleet



Efficient Fleet



Reliability

Services - Respond with low latency

Reliability

Services - Respond with low latency

Batches - Schedule with low latency

Reliability

Services - Respond with low latency

Batches - Schedule with low latency

Failures are rare in frequency (TBF)

Reliability

Services - Respond with low latency

Batches - Schedule with low latency

Failures are rare in frequency (TBF)

Failures recover quickly (TTR)

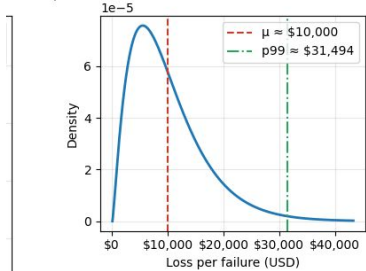
Reliability

Services - Respond with low latency

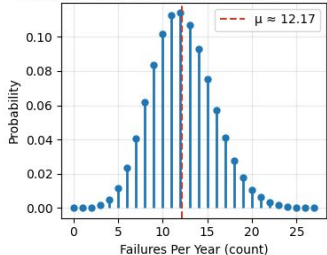
Batches - Schedule with low latency

Failures are rare in frequency (TBF)
Failures recover quickly (TTR)
Failures have low impact (Blast Radius)

Loss | Failure: Beta($\mu \approx \$10,000$, $p99 \approx \$31,494$)



Failures Per Year: Poisson($\lambda T \approx 12.17$)



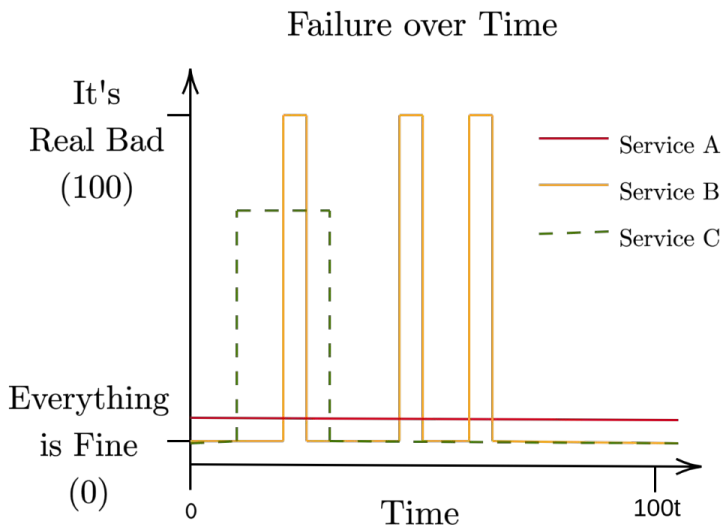
Instead of Measuring Nines, ask

How often does it fail (TBF)?

What is the impact (Blast Radius)?

How long to recover (TTR)?

Nines are Deceptive



Different Failure Modes Require Different Solutions!

	Service A	Service B	Service C
Nines	99.999	99.999	99.999
Time Between Failure ($P_{05} + P_{50}$)	(0, 0)	(3t, 10t)	(10t, 40t)
Time to Recovery (P_{05}, P_{50})	(∞ , ∞)	(1t, 2t)	(4t, 5t)
Impact of Failure (P_{05}, P_{50})	(1, 1)	(100, 100)	(50, 80)
Potential Resolution	Hedging or Tests	Load Shedding or Backpressure	Alerting or Failover or Bl...



N

Understand
Hardware Supply

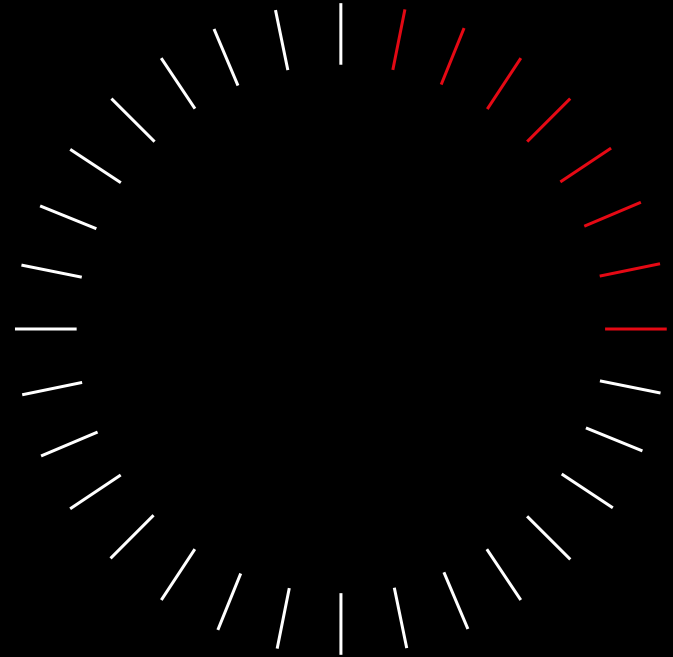
Understand
Software Demand

Balance
Compute Supply
Software Demand

Understand Supply

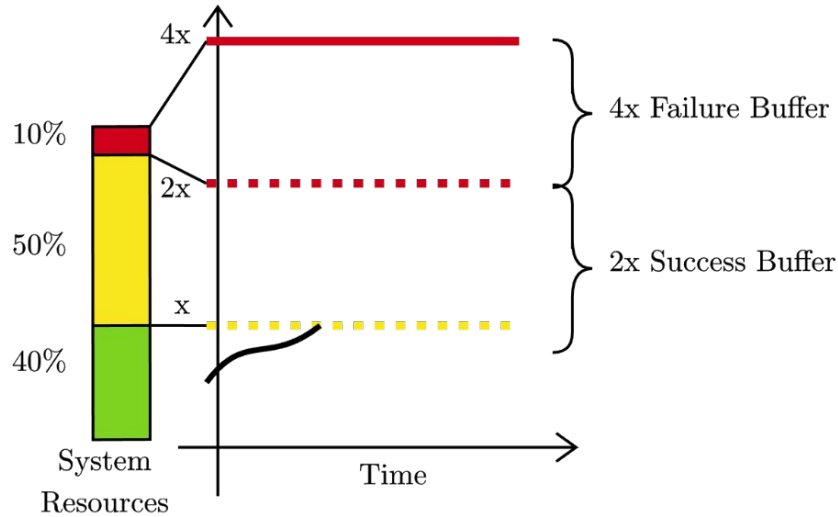
Capacity Planning

Hardware Supply Signals

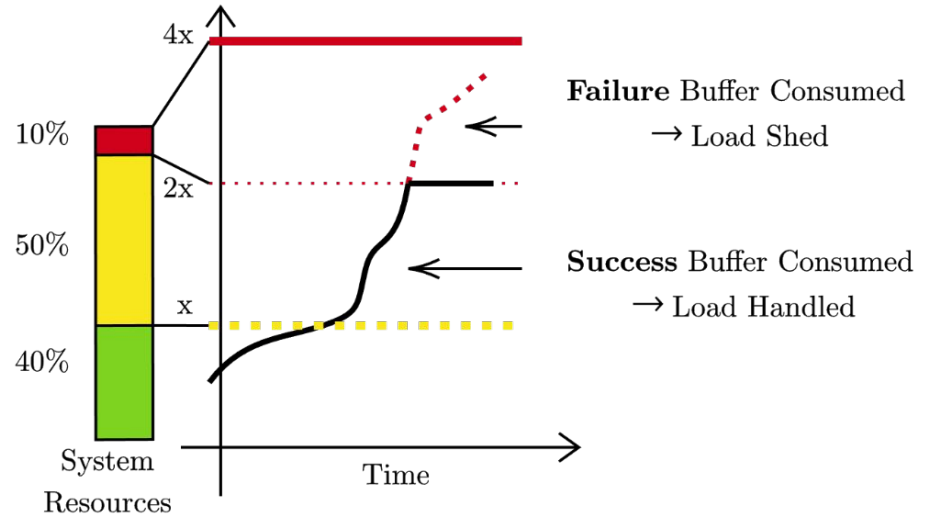


Reason about Headroom with **Buffers**

Normal System Load with Buffer

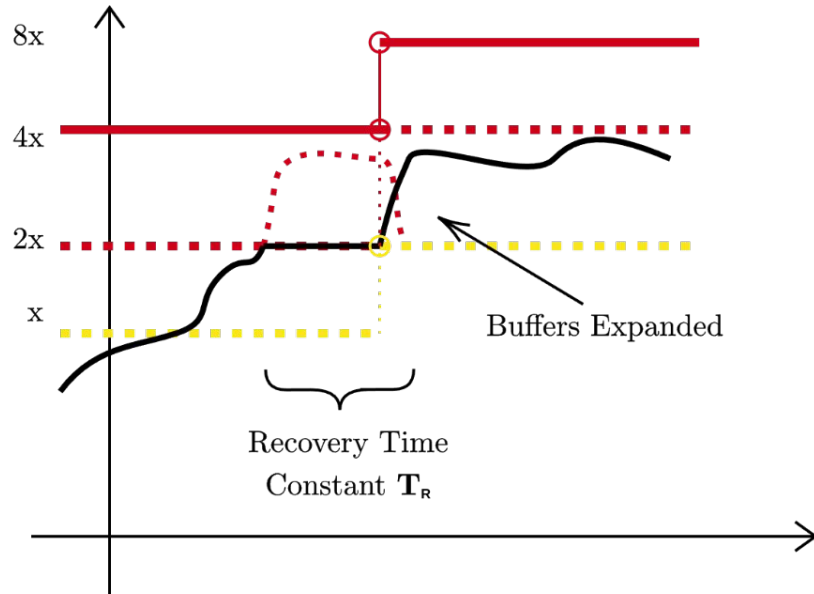


System Load Under Load Spike



Buffer is Linked to Business Outcomes

Buffer Recovering After Load Spike

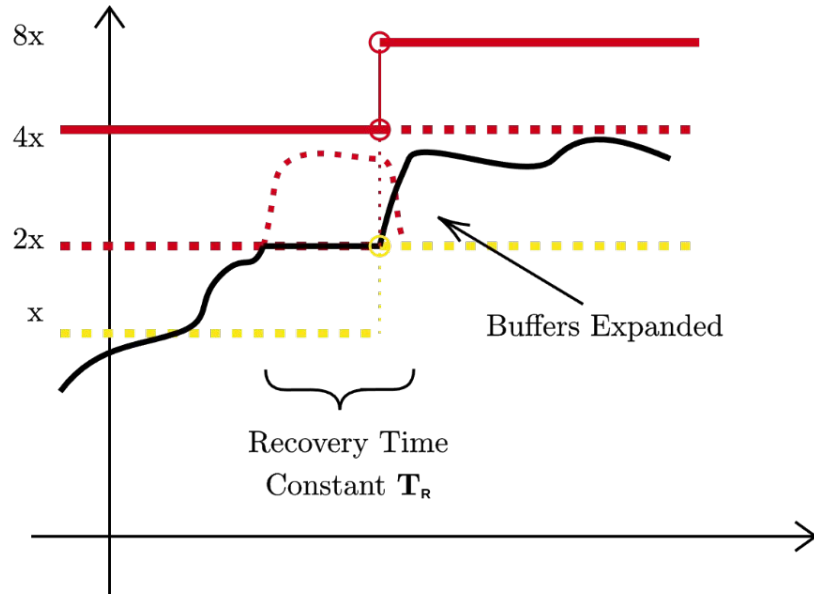


Buffer = F(
Criticality Tier,
Business Domain,
Recovery Time Constant
)

Low utilization is a tradeoff!

Buffer is Linked to Business Outcomes

Buffer Recovering After Load Spike



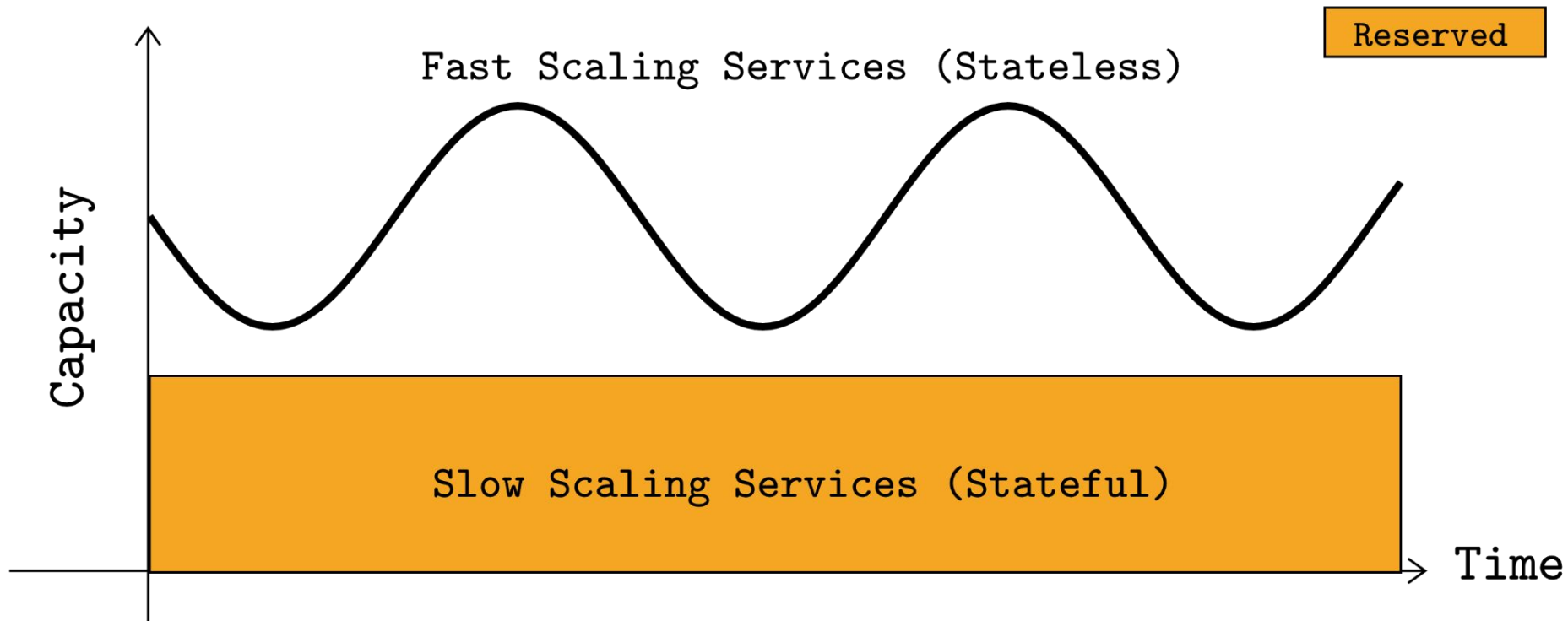
Buffer = F(
Criticality Tier,
Business Domain,
Recovery Time Constant
)

Low utilization is a tradeoff!

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

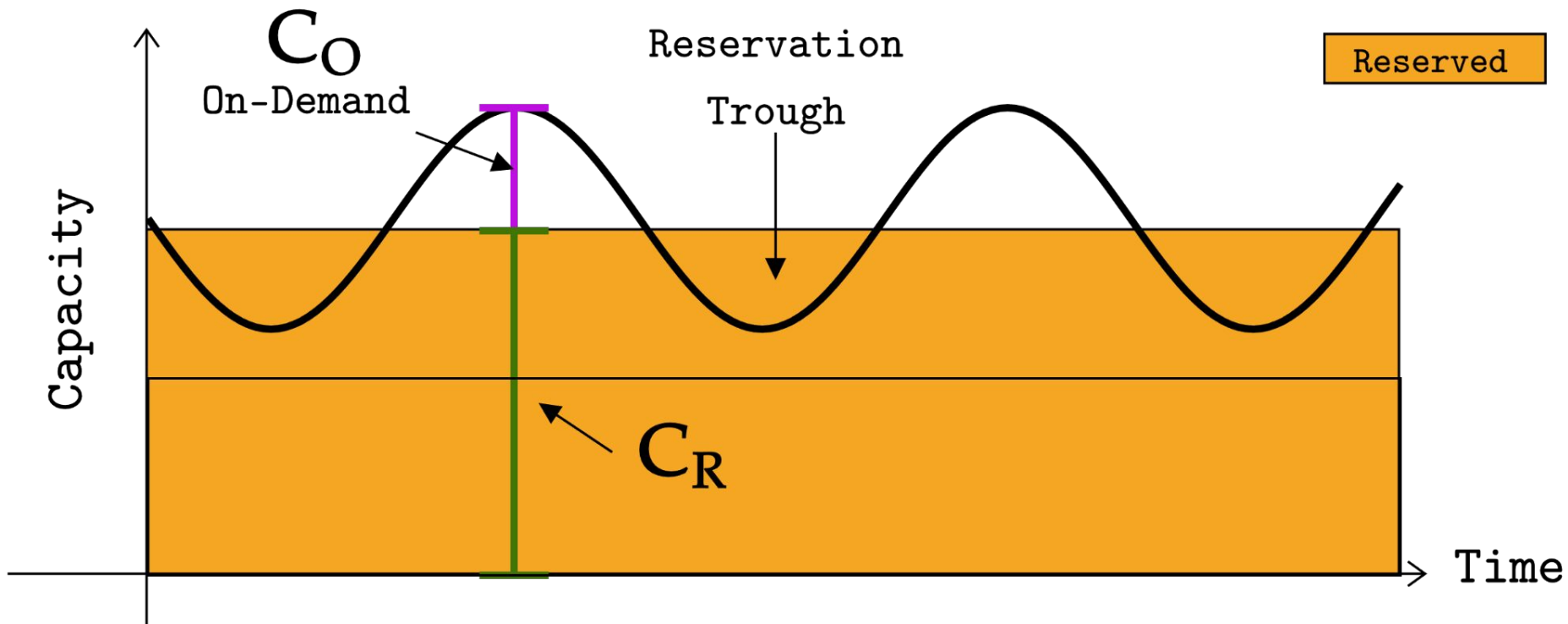
Cloud Reality #1

Still have to Plan



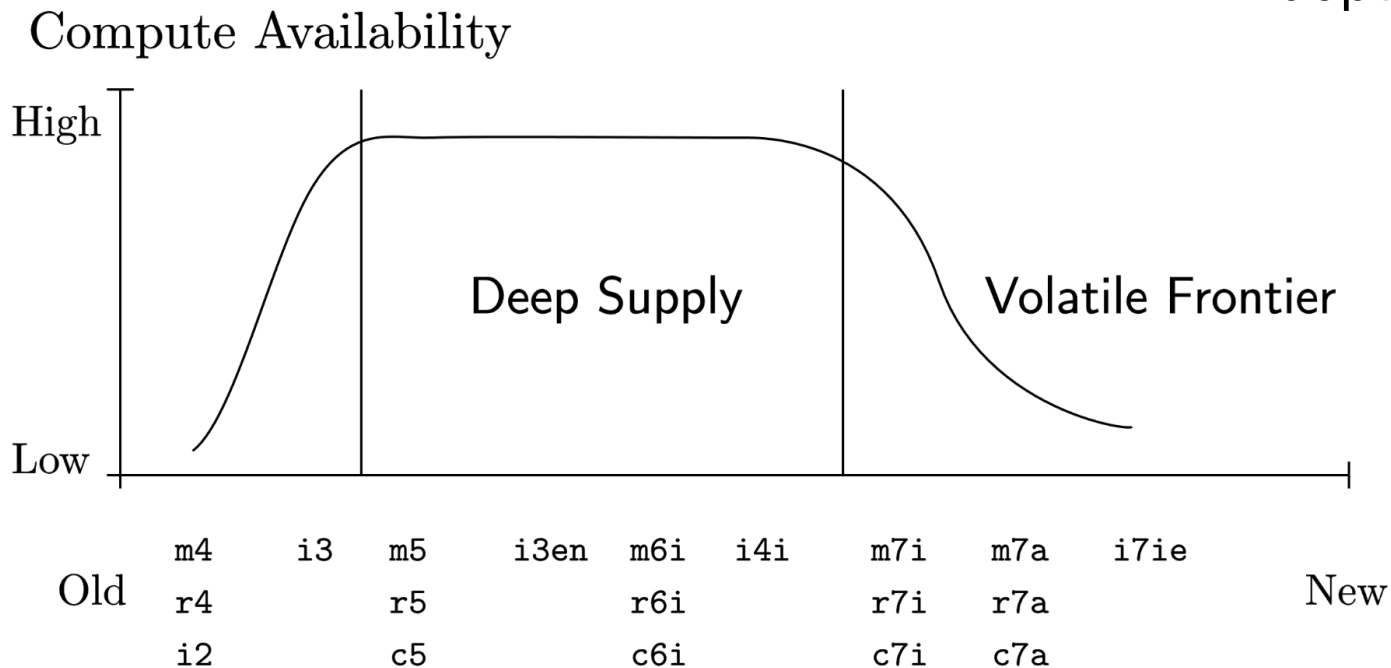
Cloud Reality #1

Still have to Plan



Cloud Reality #2 Computers have *Variable Supply*

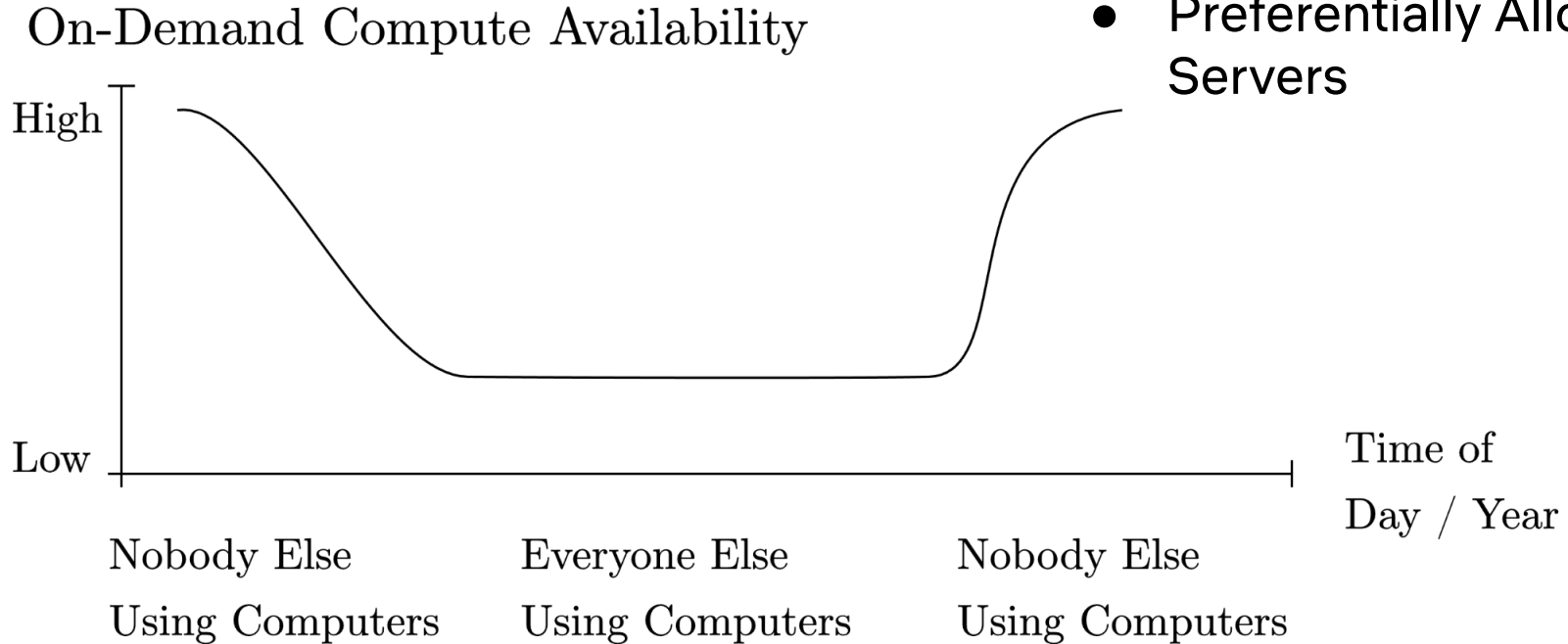
Have to *carefully*
Adopt New Shapes



Cloud Reality #2 Computers have *Variable Supply*

Options:

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

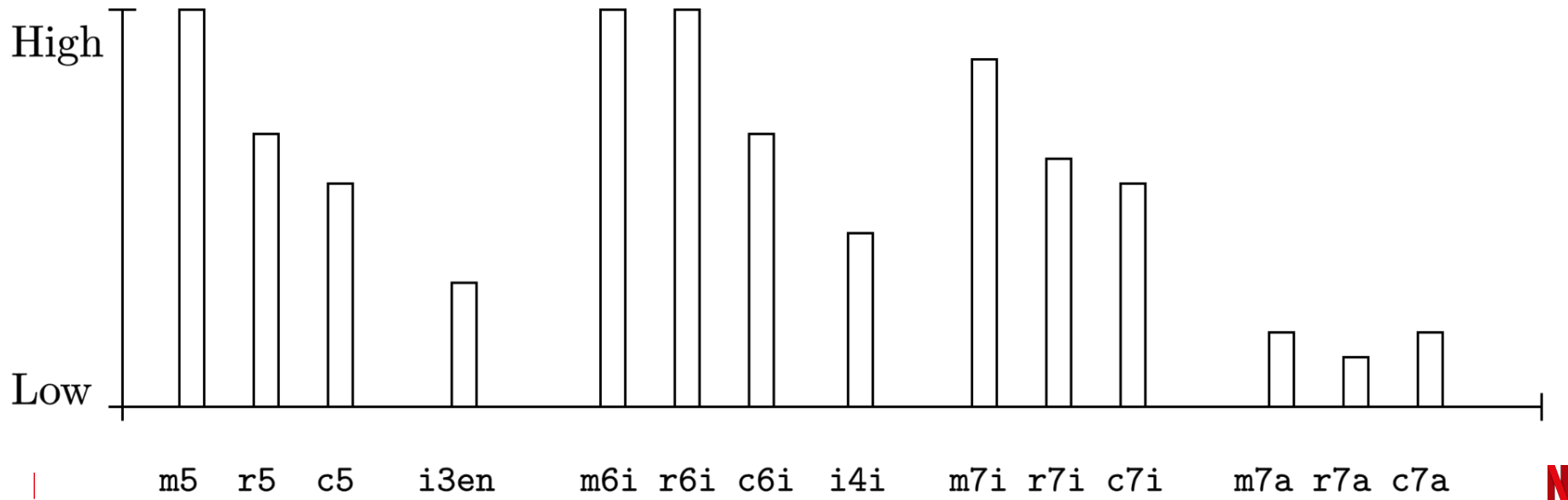


Cloud Reality #2
Computers have
Variable Supply

More flexibility == More Capacity

Assign workloads by resources not names

Compute Availability by Shape - Data Not Real

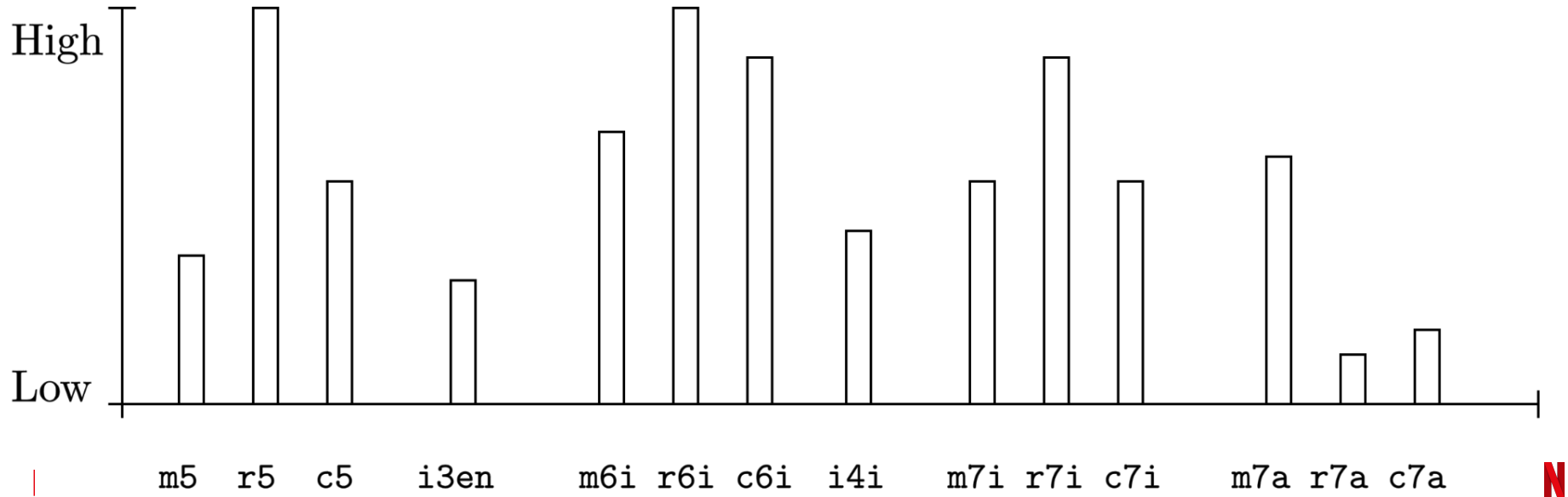


Cloud Reality #2
Computers have
Variable Supply

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



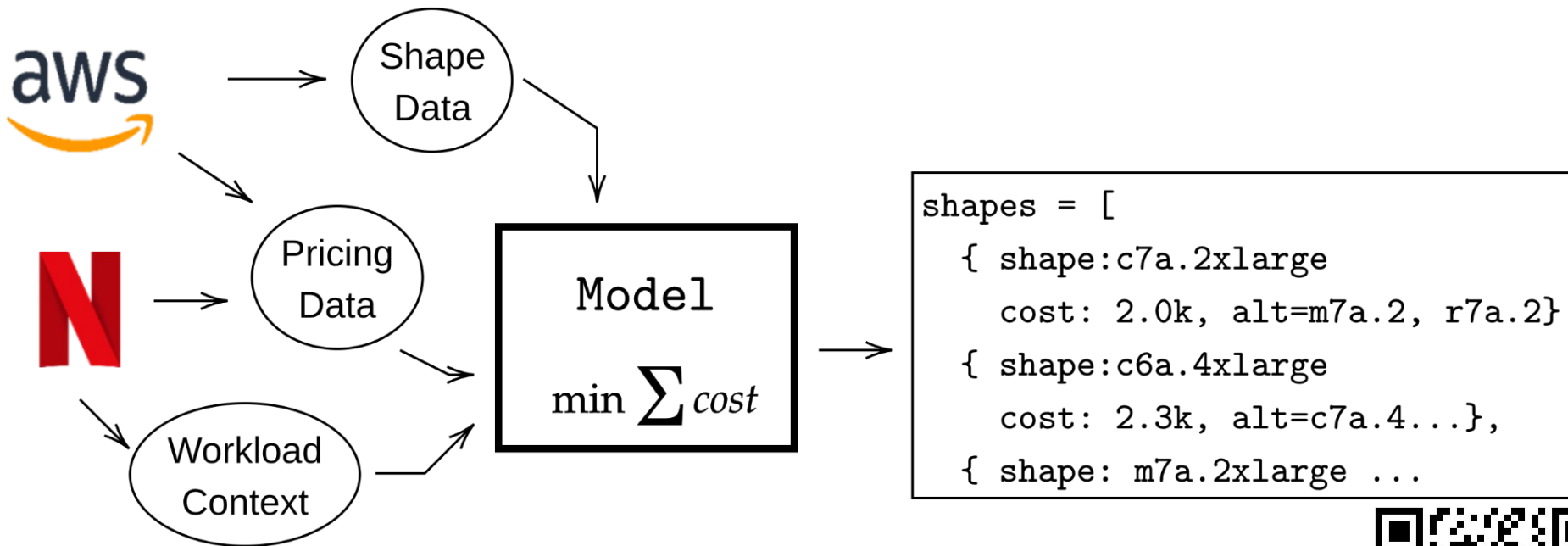
Service Capacity Modeling

Build passing

A generic toolkit for modeling capacity requirements in the cloud. |

```
"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
    }
  }
},
},
```

Capacity Model



Capacity Plan

Strategy: Buffer, and lots of math

```
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
)
```



```
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}
```

Buffer Differs By Server Type

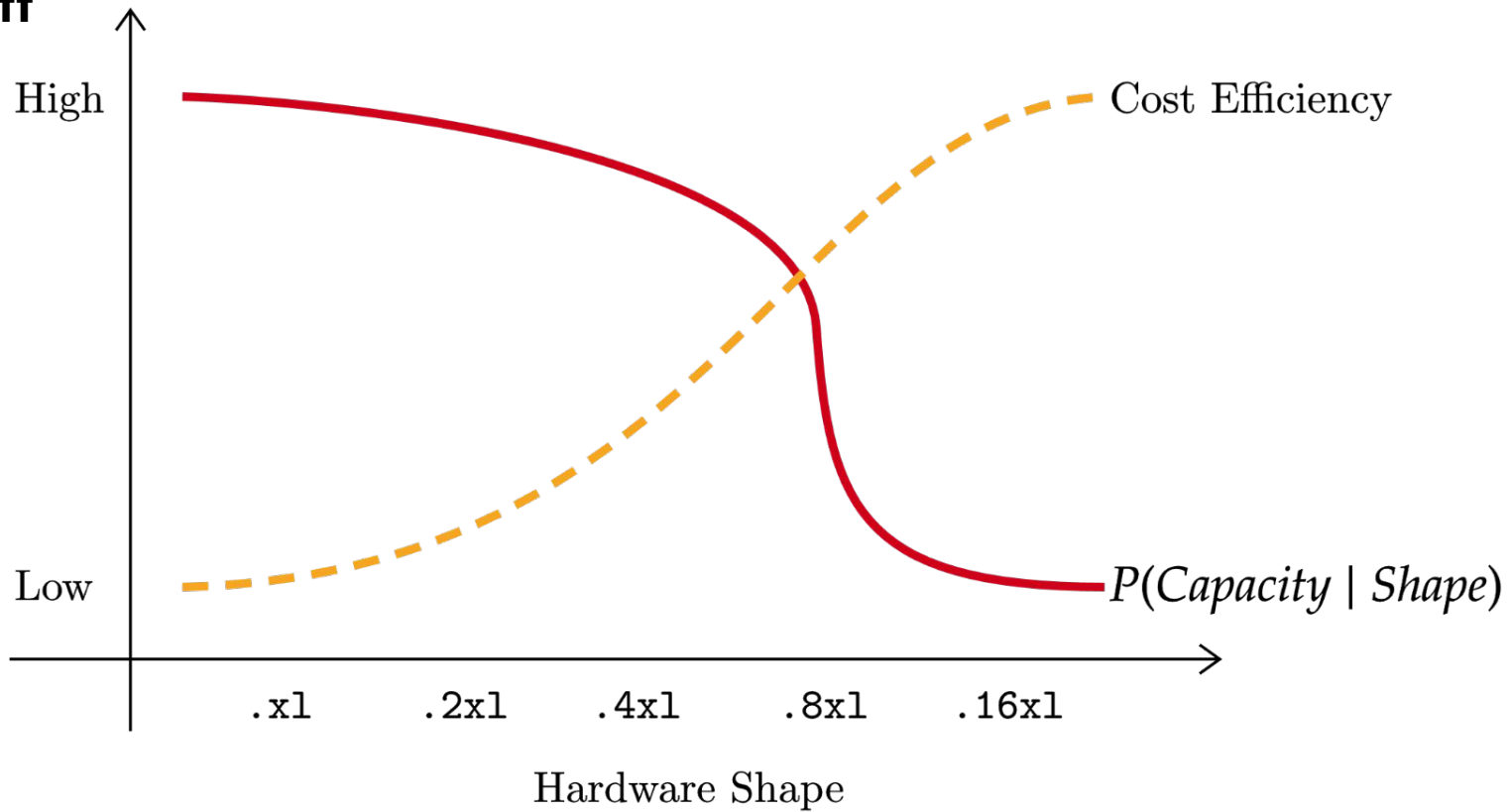
	Name	Failure Buffer %	Tier 1 Target %	Tier 0 Target %
19	m7a.xlarge	0.70	0.46	0.35
21	m7a.4xlarge	0.84	0.56	0.42
22	m7a.8xlarge	0.88	0.59	0.44
29	m7i.xlarge	0.62	0.41	0.31
31	m7i.4xlarge	0.79	0.53	0.40
32	m7i.8xlarge	0.85	0.57	0.42

Buffer Grid

CPU Targets: Instance Types vs Buffers (SVOD, TIER_0)

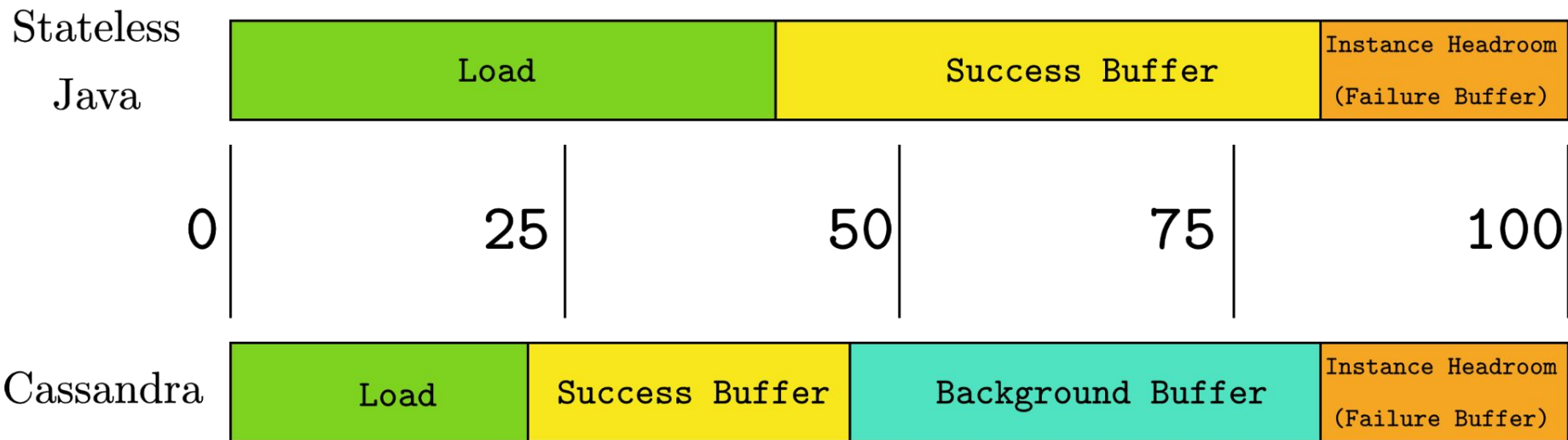


Capacity Tradeoff



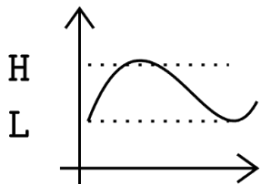
Buffer Differs By Workload

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



Ensure Supply

Long
Running
Variable



Financial

$$\min \sum \mathbb{E}[cost]$$
$$\mathbb{E}[cost] \propto \$ \times T$$

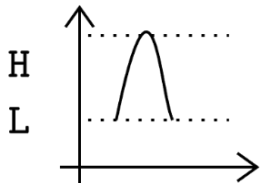
Long
Running
Stable



Capacity

$$\min \sum \mathbb{E}[risk]$$
$$\mathbb{E}[risk] \propto \frac{Loss}{P(Capacity)}$$

Oneshot



Reserve (

[L, H]:m7a.4xlarge
) → Launch

Reserve (

N:m7a.4xlarge →

Capacity Reserve (

[L, H]:m7a.4xlarge
) → On-Demand



Takeaways

Buffer is the relevant configuration

Capacity must be prioritized

Time-to-Recover matters

Compute is Fungible via the **Right Math™**

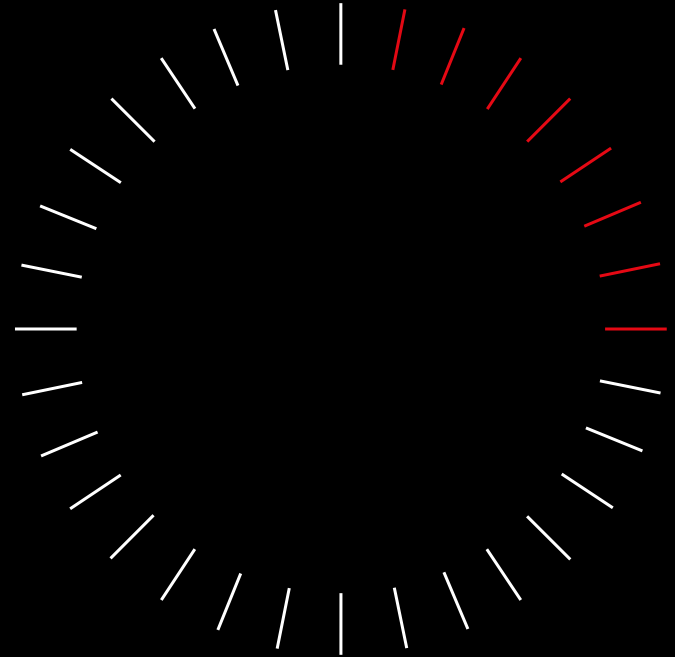
Understand
Hardware Supply

Understand
Software Demand

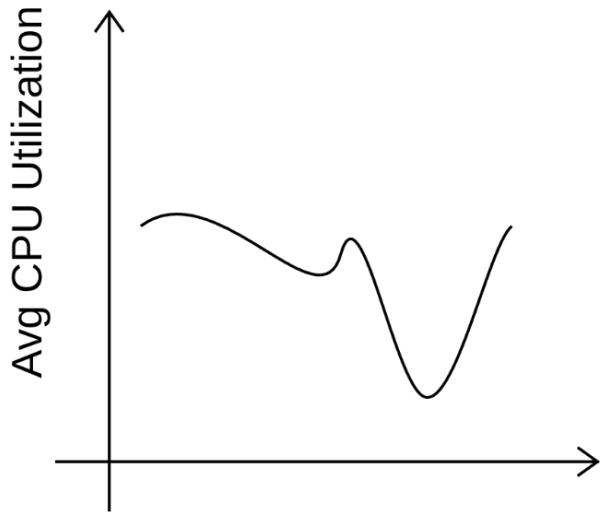
Balance
Compute Supply
Software Demand

Understand Demand

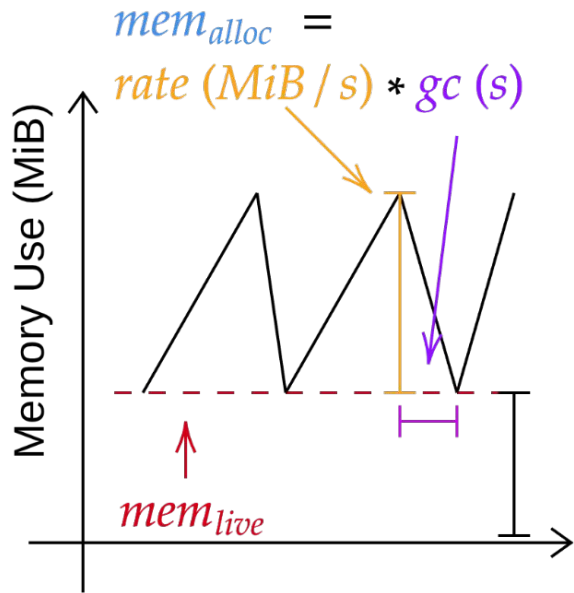
Workload Profile



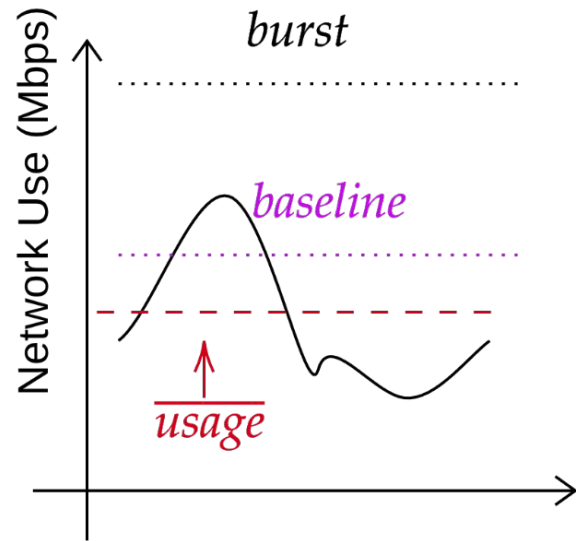
Workload analysis



$$CPU_{util} = \max(\overline{cpu_{util}})$$

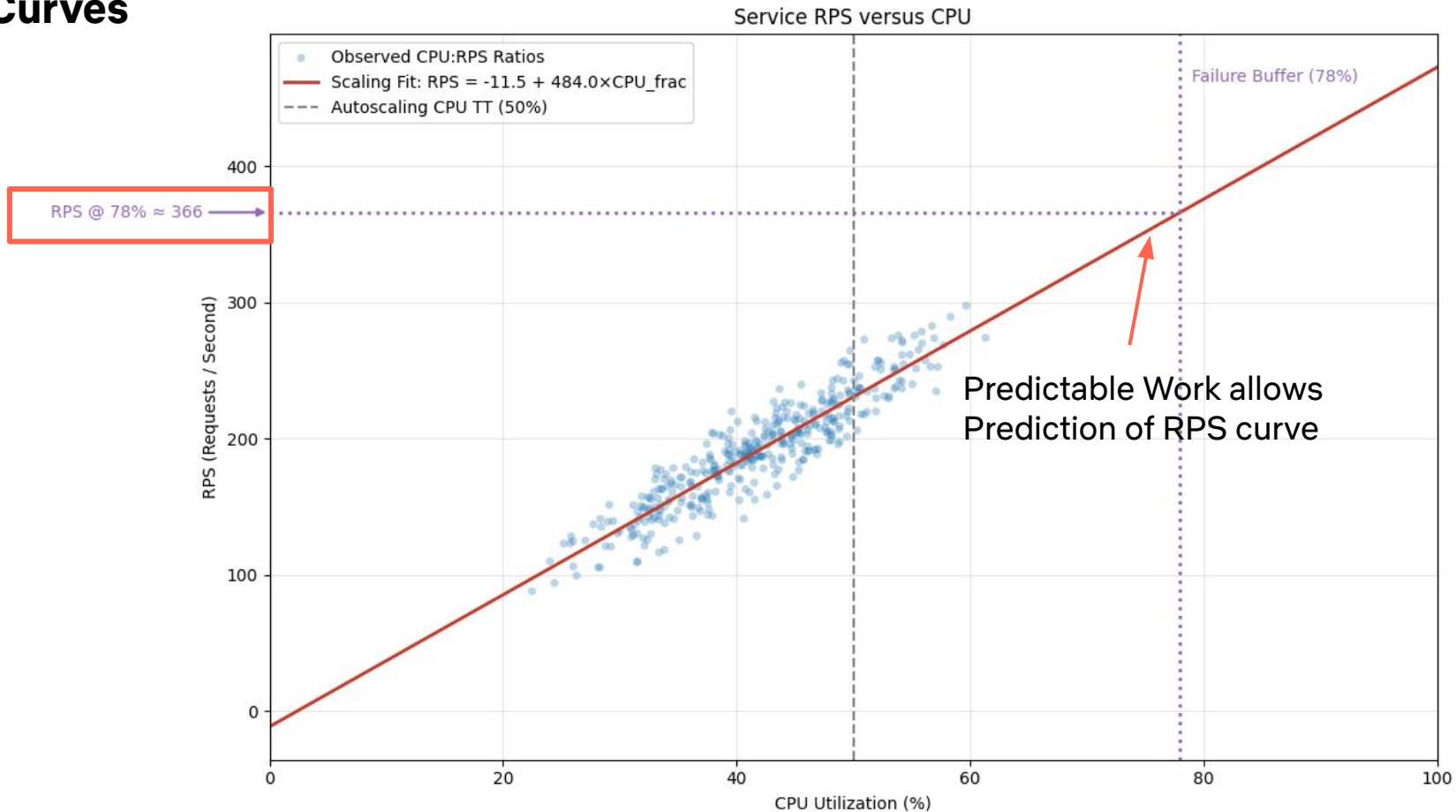


$$Mem = mem_{live} + mem_{alloc}$$

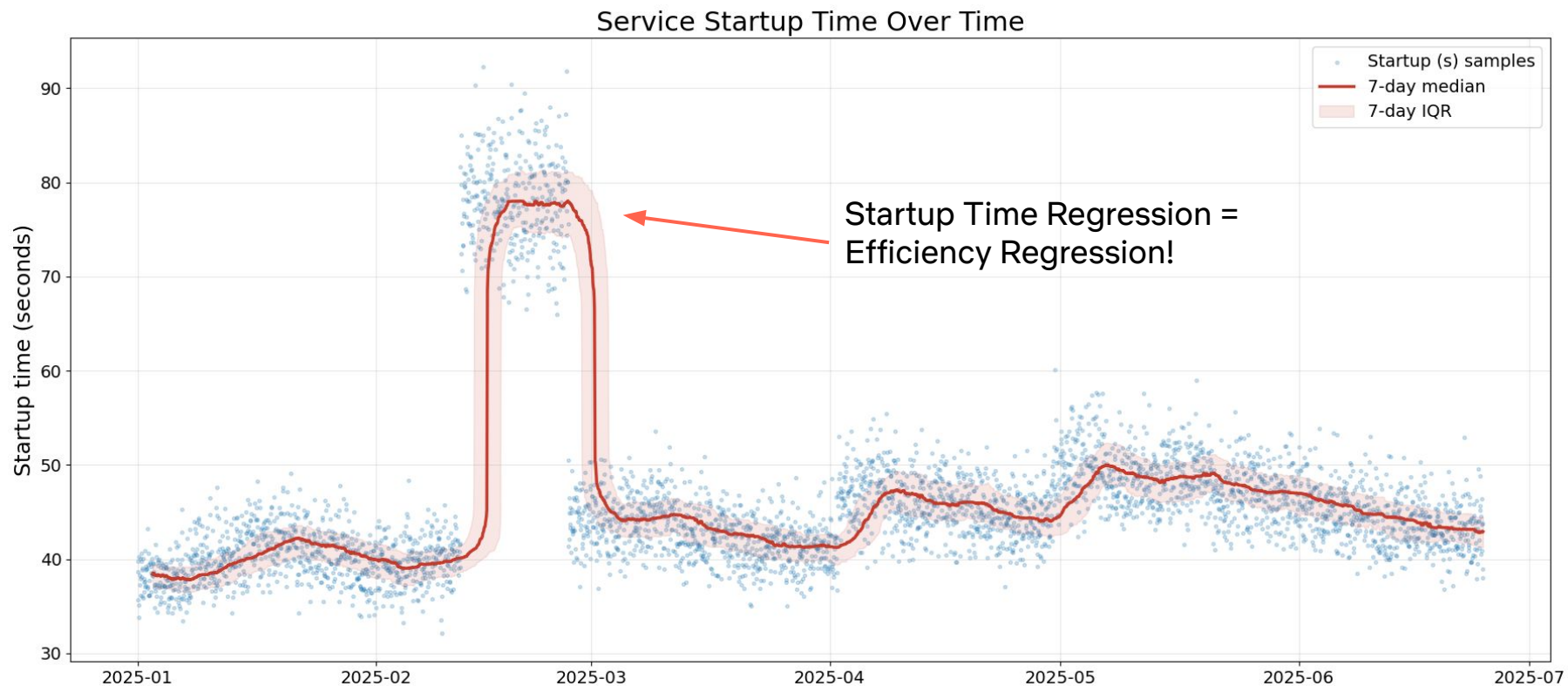


$$Net_{util} = \frac{\overline{usage}}{baseline}$$

Scaling Curves



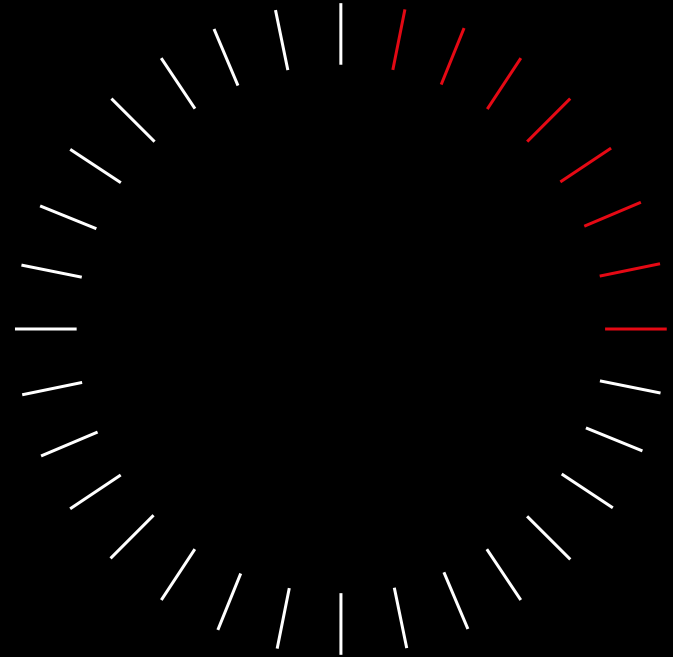
Startup Times



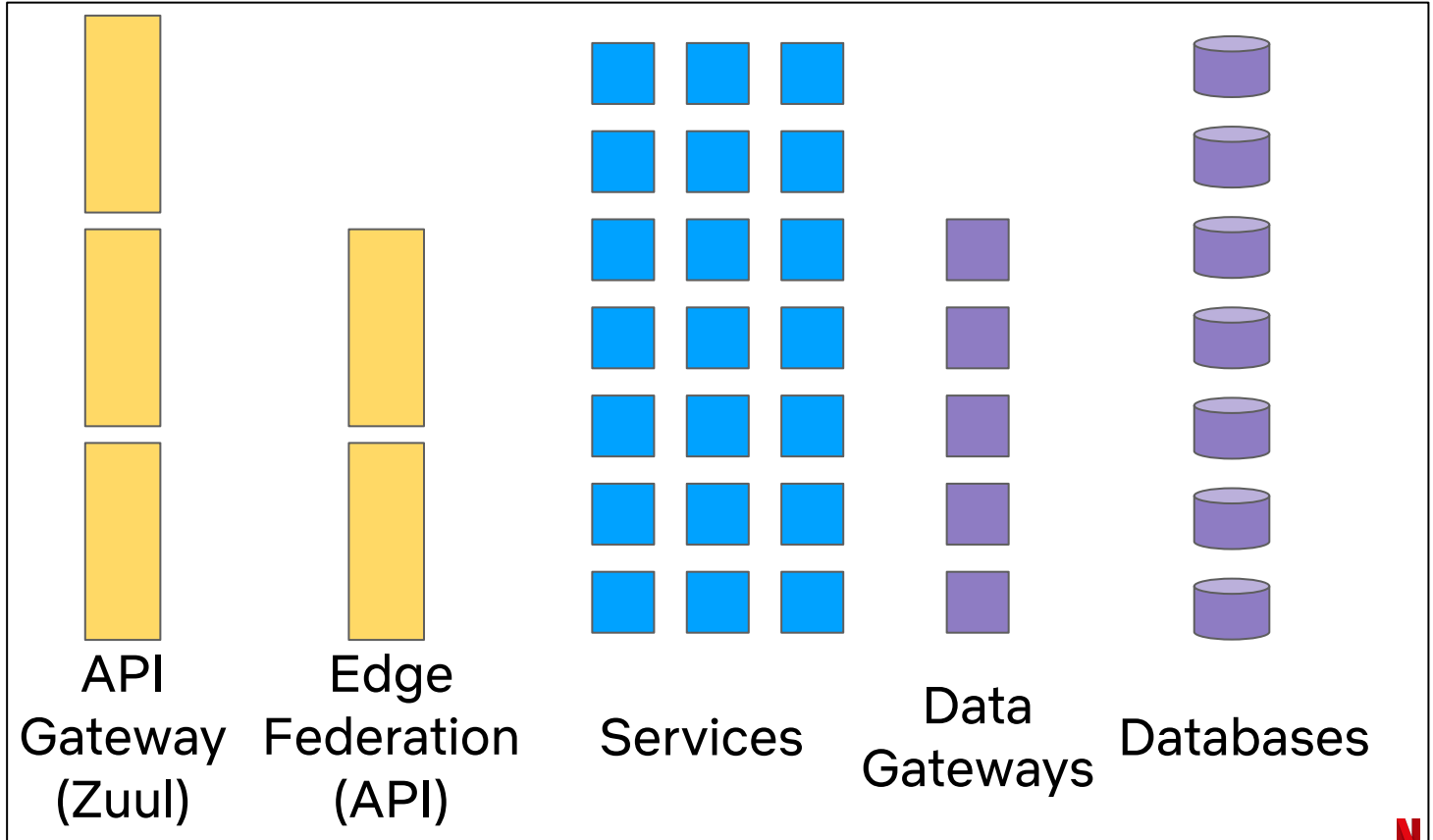
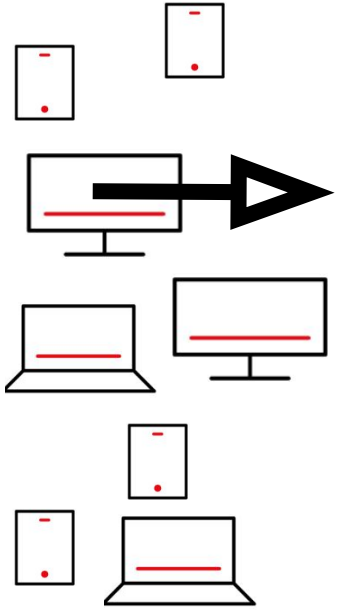
Understand Demand

Workload Profile ✓

Traffic Analysis

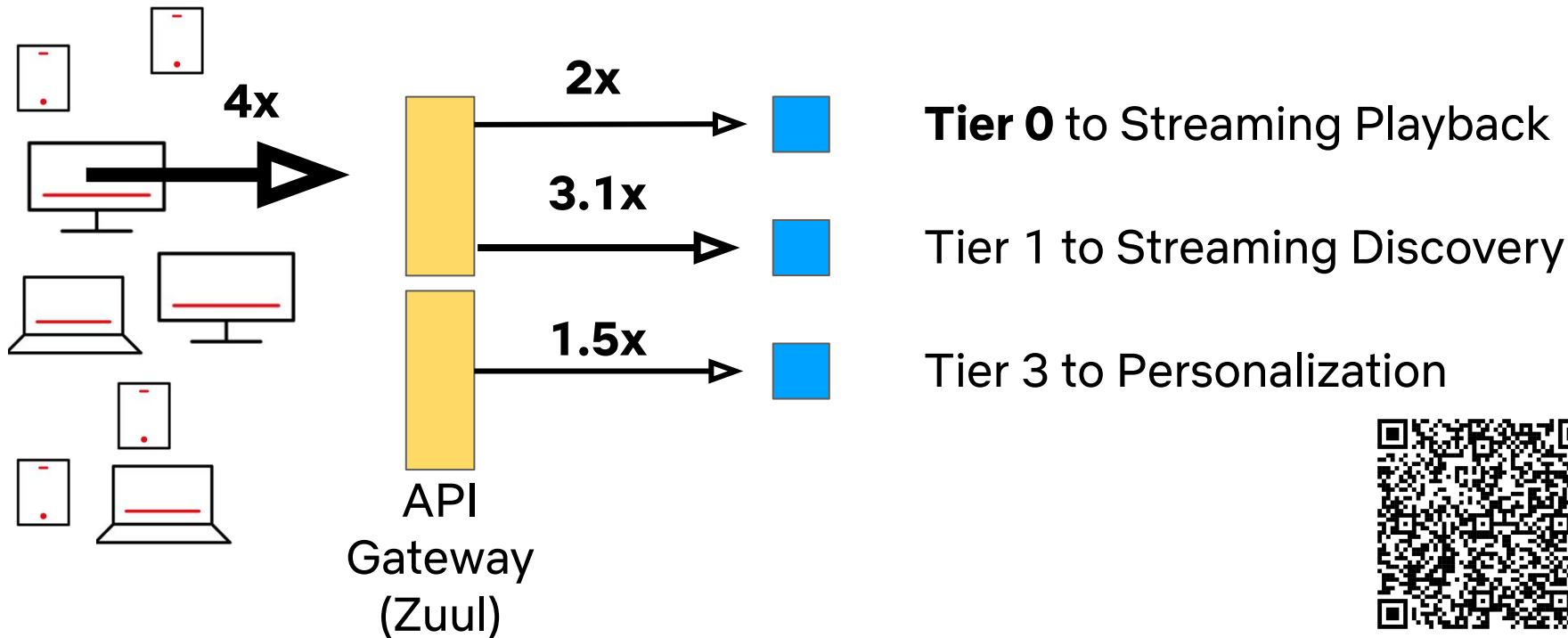


Microservices

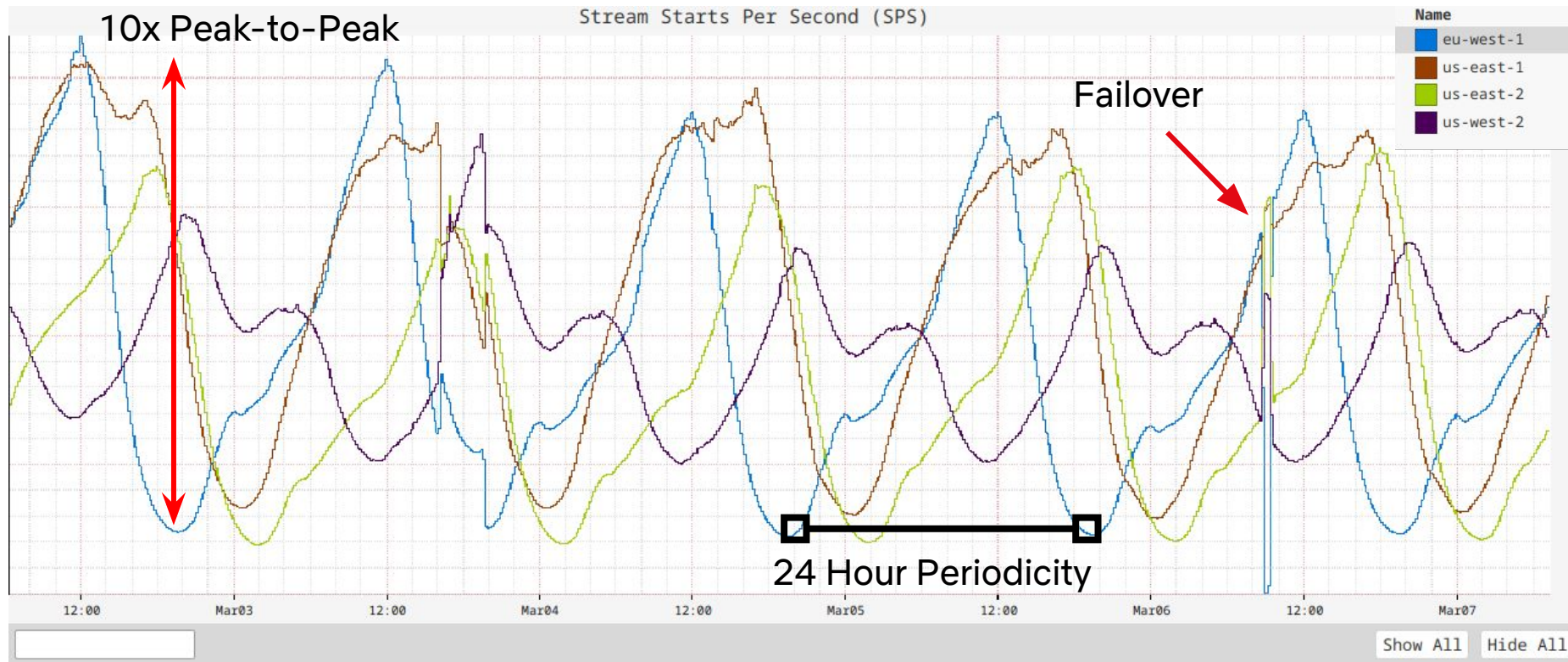


Different for Every Service

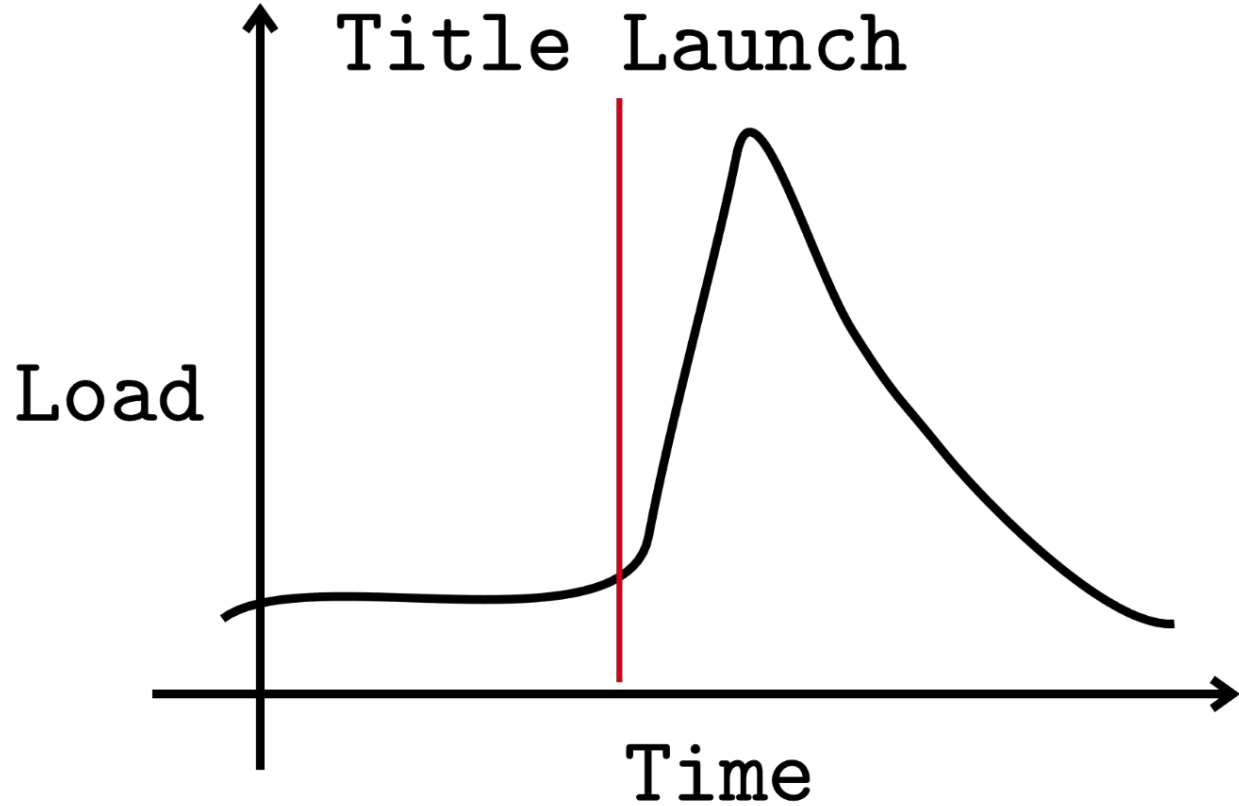
Call Graph is not Uniform in Traffic or Criticality



Variable Starts-Per-Second (SPS) Load



Content Driven Demand



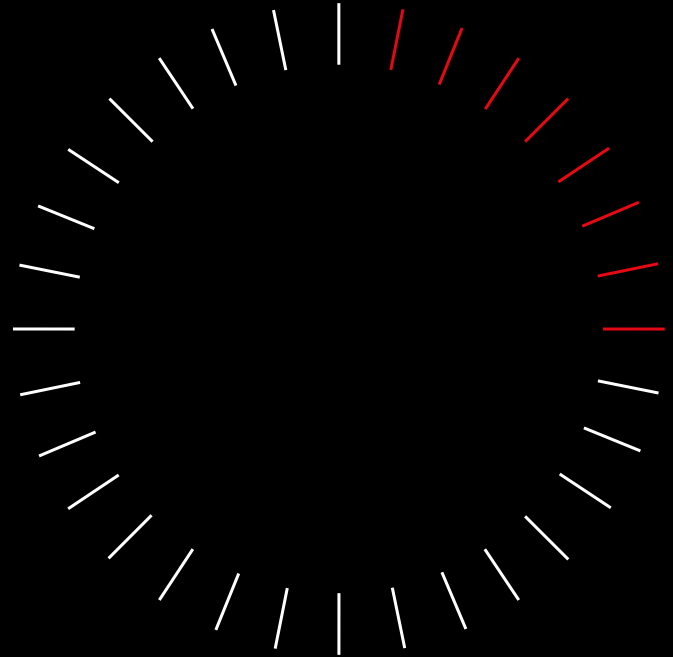
Understand
Hardware Supply

Understand
Software Demand


Balance
Manage Supply
Manage Demand

Tenets

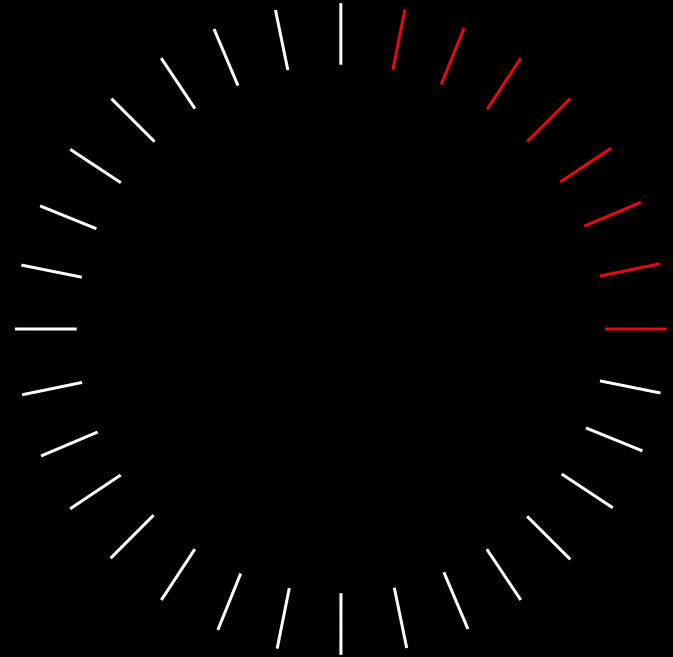
Global vs local optimization



Tenets

Global vs local optimization 

Ensure Buffers exist for **critical** services, when **it matters**

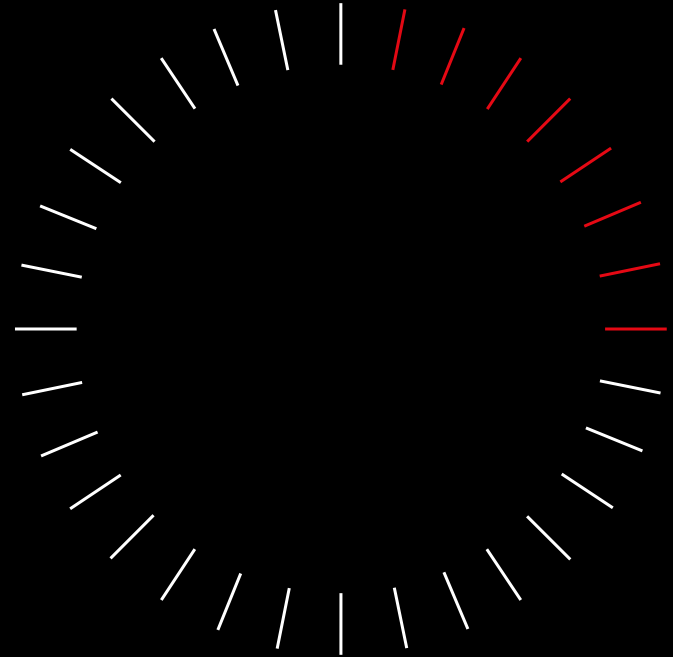


Tenets

Global vs local optimization ✓

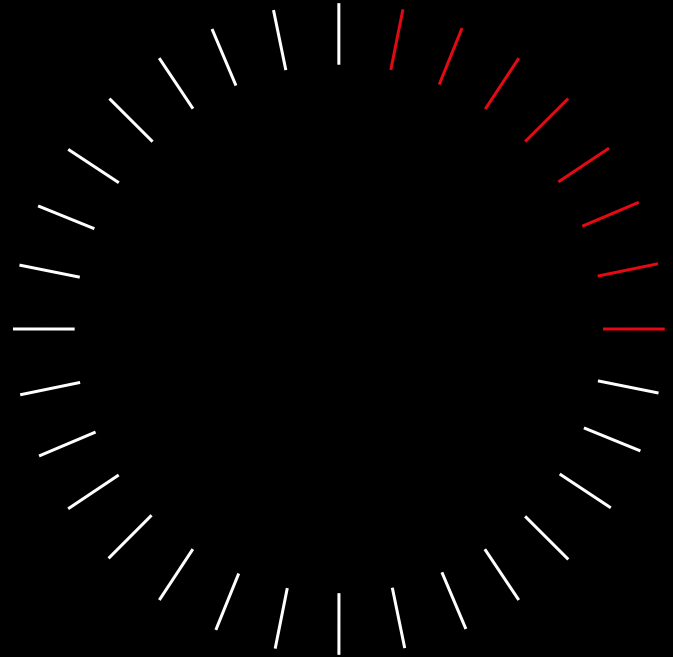
Buffers exist for **critical** services,
when **it matters** ✓

Managed Supply and Demand

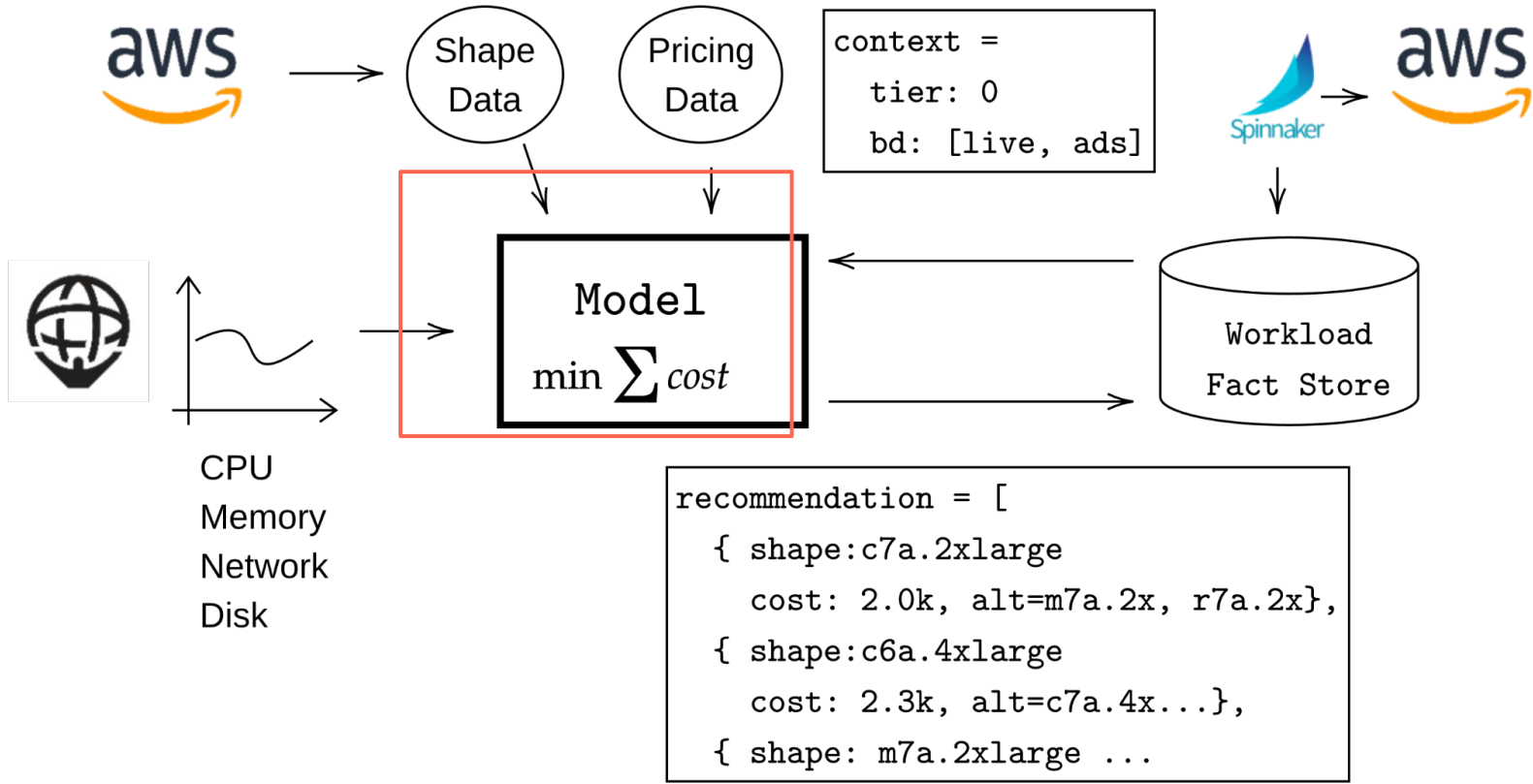


Manage Supply

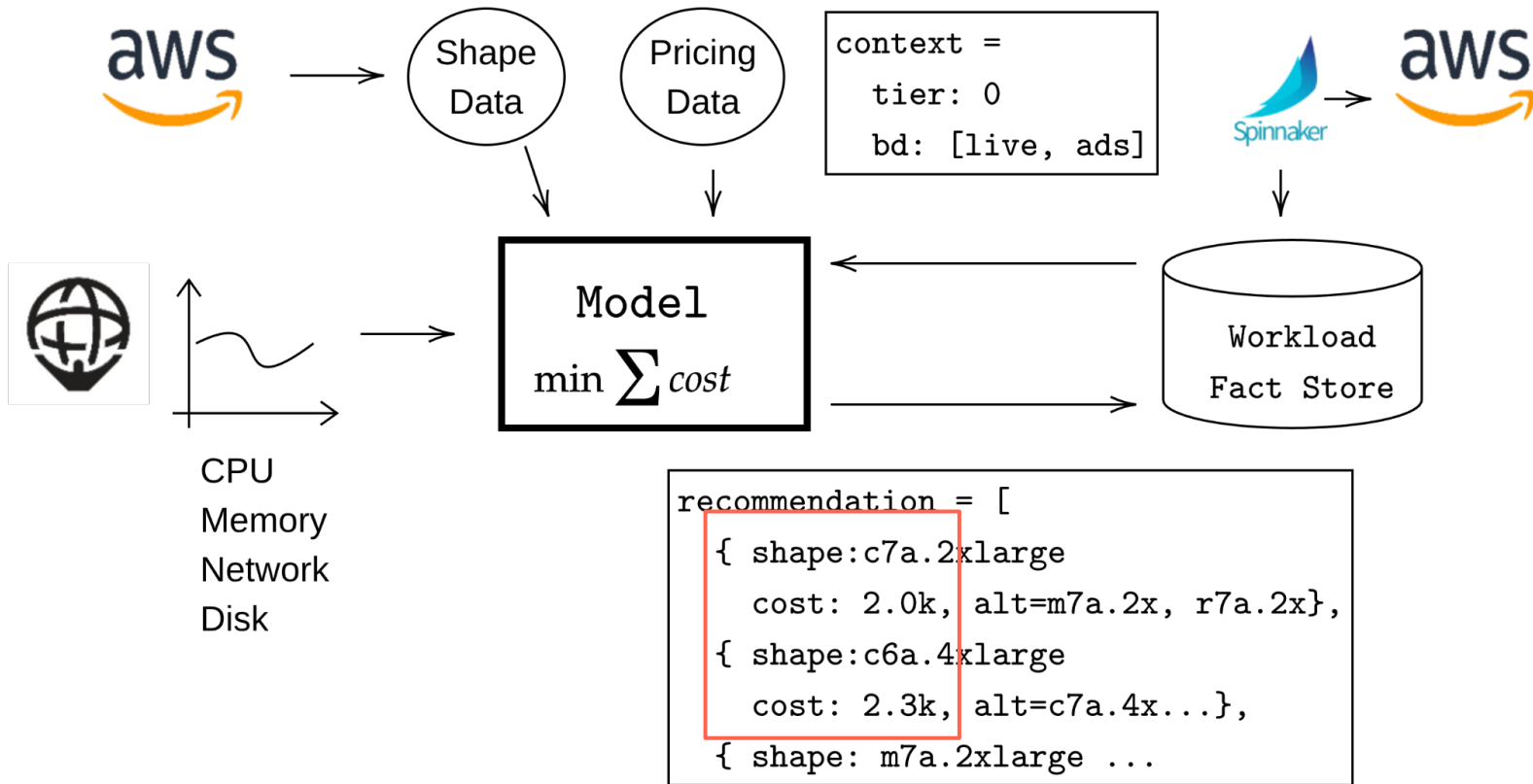
Hardware 🤝 Workload



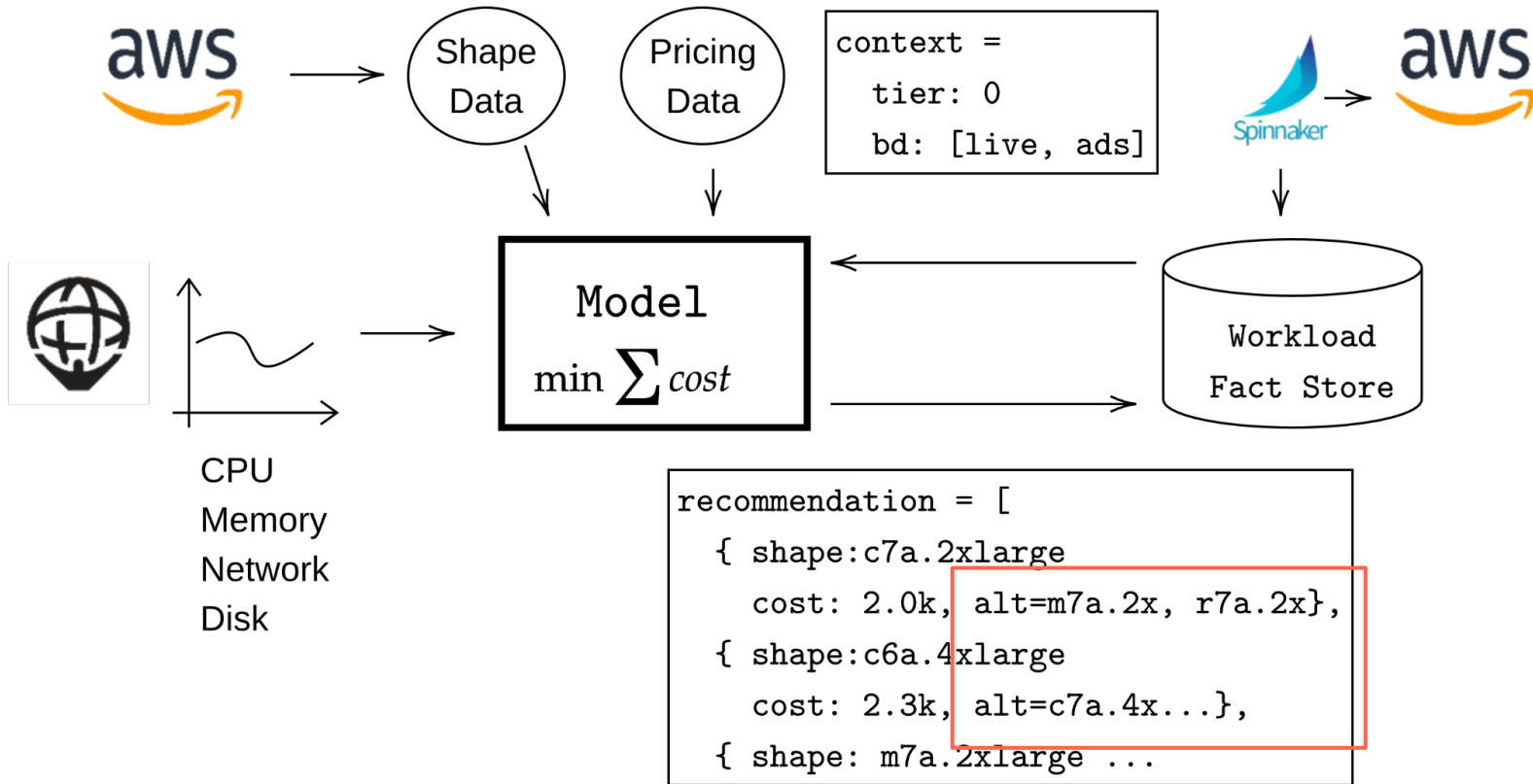
Hardware Selection



Hardware Selection



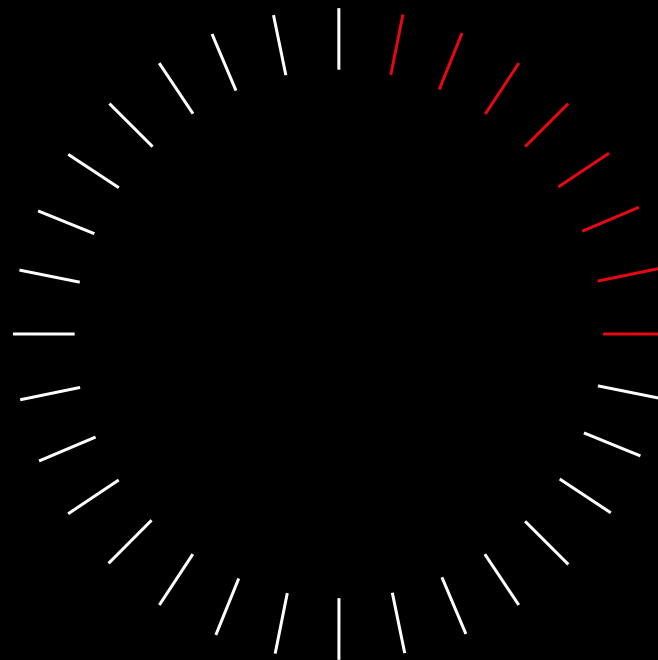
Hardware Selection



Manage Supply

Hardware 🤝 Workload ✅

Fleet **Shaping**



Fleet Shaping

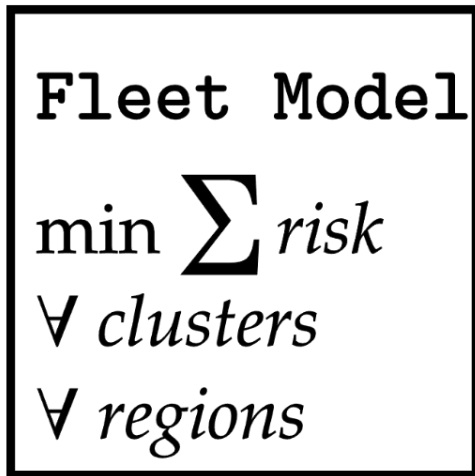
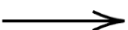
Perf/Cost
Context



Capacity
Context



Business
Context

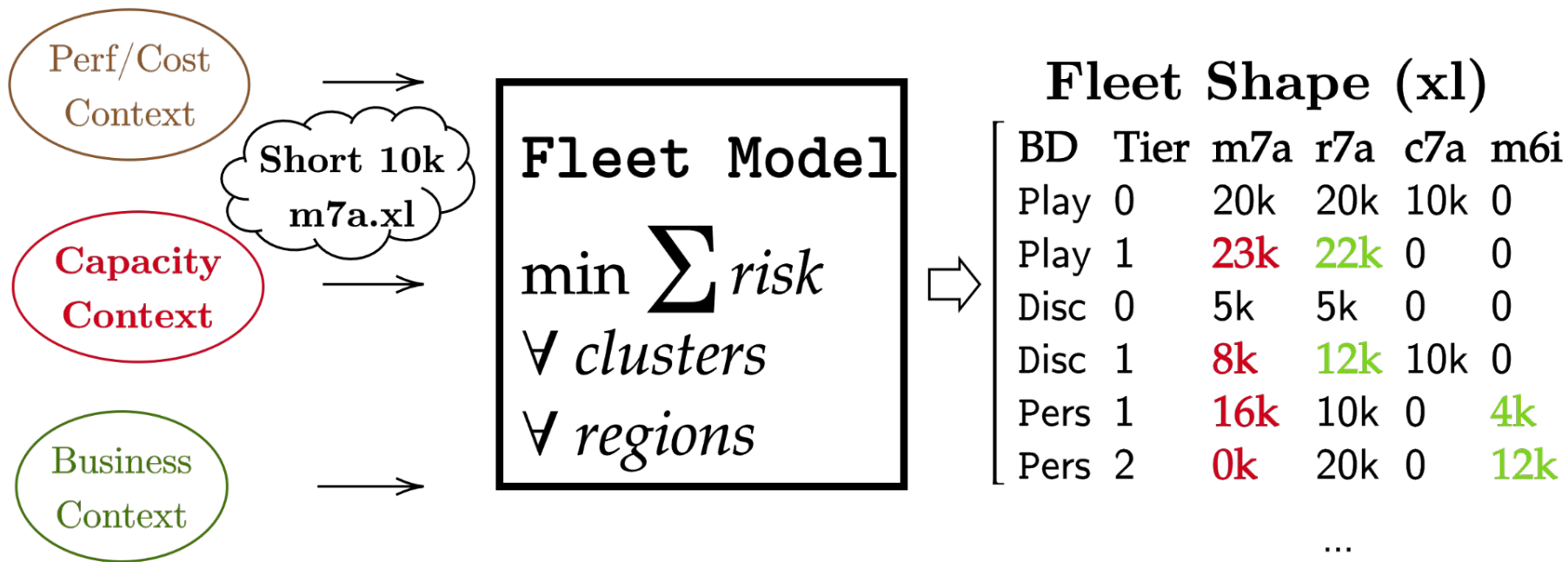


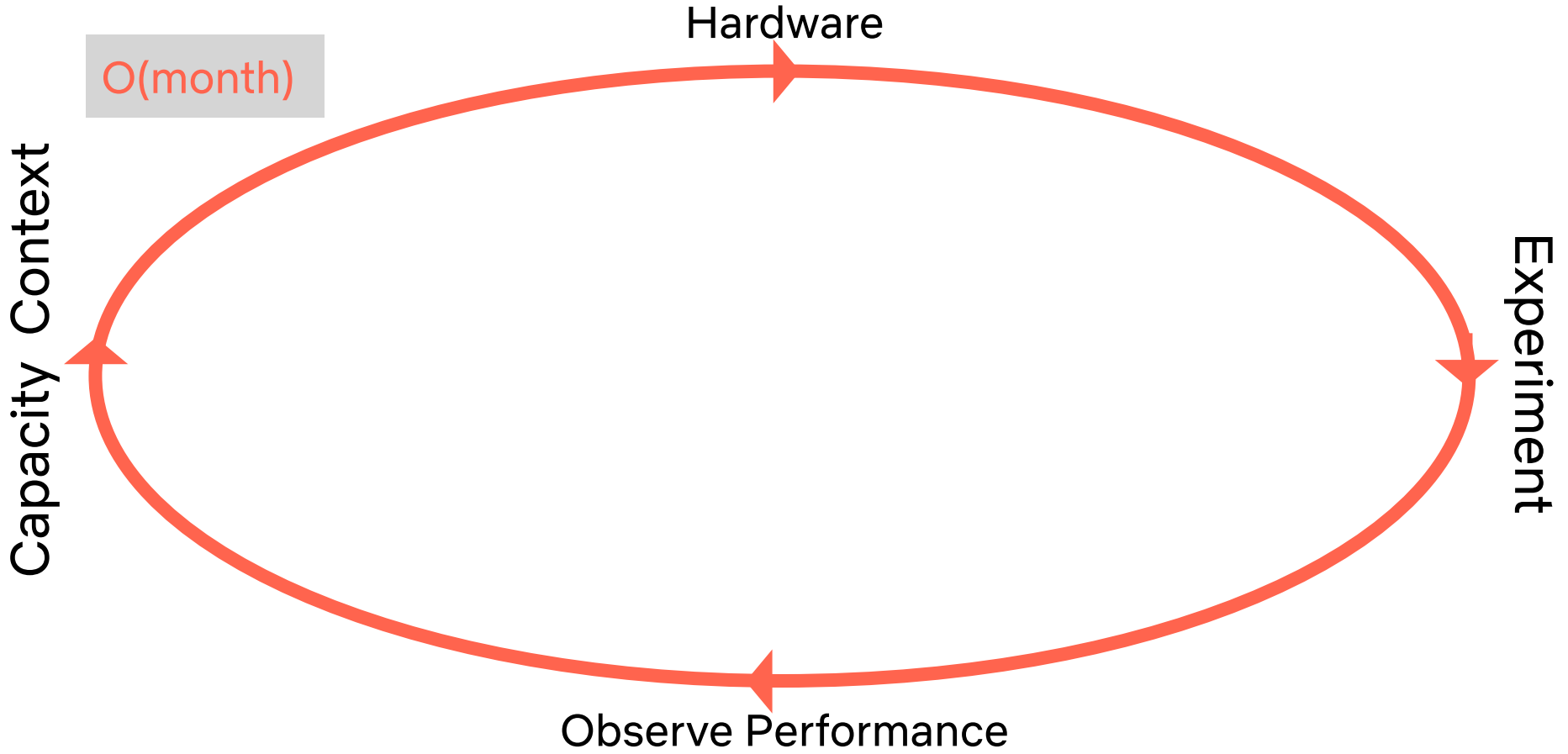
Fleet Shape (x1)

BD	Tier	m7a	r7a	c7a	m6i
Play 0		20k	20k	10k	0
Play 1		25k	20k	0	0
Disc 0		5k	5k	0	0
Disc 1		10k	10k	10k	0
Pers 1		20k	10k	0	0
Pers 2		2k	20k	0	10k

...

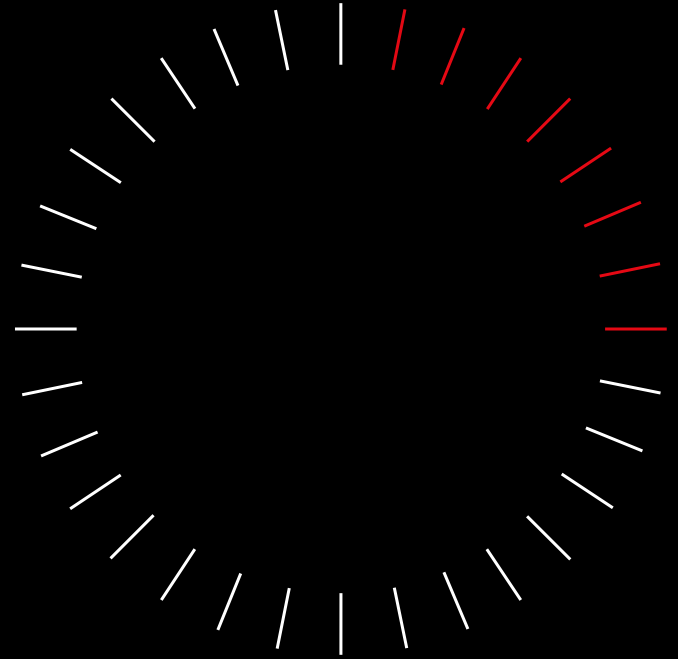
Fleet Shaping





Manage Demand

Pre-scaling

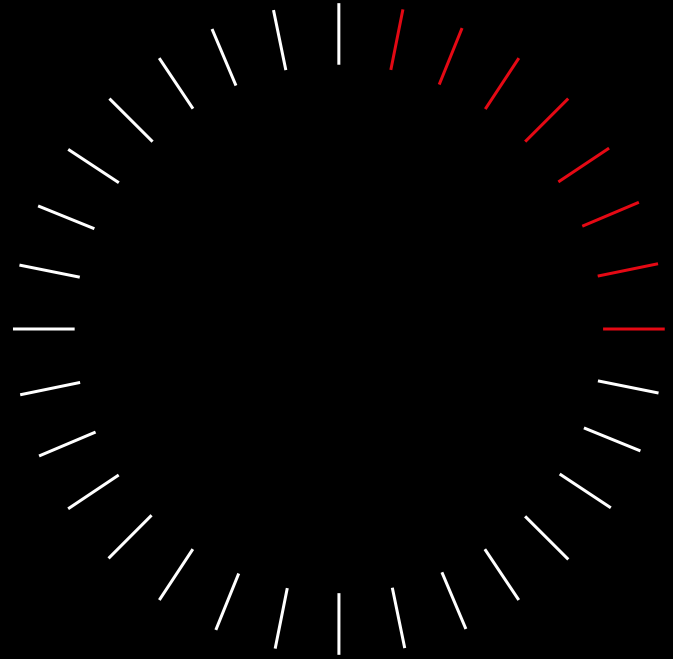


Pre-scaling

Prediction models 🧙‍♂️

Peak load estimation (viewers -> traffic)

Automated fleet scaleup 🤖



Pre-Scale for Load



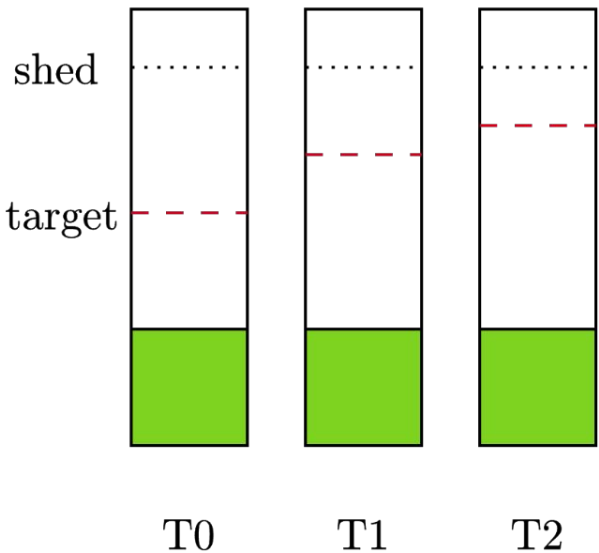
1. Pin Mins Up

2. Traffic arrives

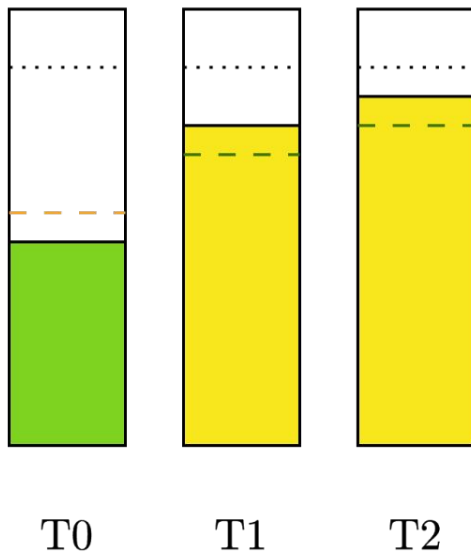
3. Unpin Mins

4. Scale Down

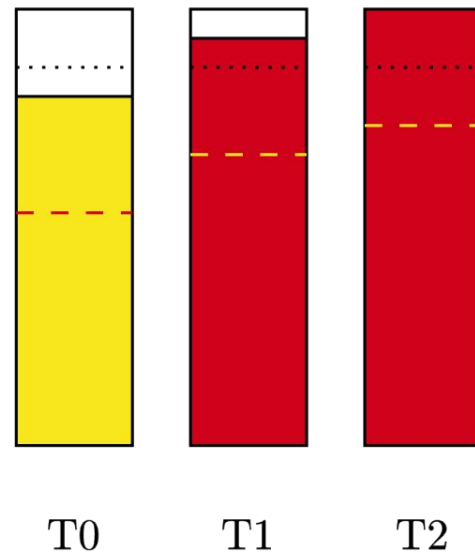
Fleet Prescale



Inefficient: Overscaled



Efficient: Balance



Probably Underscaled

Fast Recovery Mitigates Pre-Scale Risk

Load Shedding vs Autoscaling vs RPS per Instance



Axis 0

Shed RPS

Max :	1.872k	Min :	0.000
Avg :	156.689	Last :	0.000
Tot :	5.014k	Cnt :	32.000

Axis 1

Fleet Size

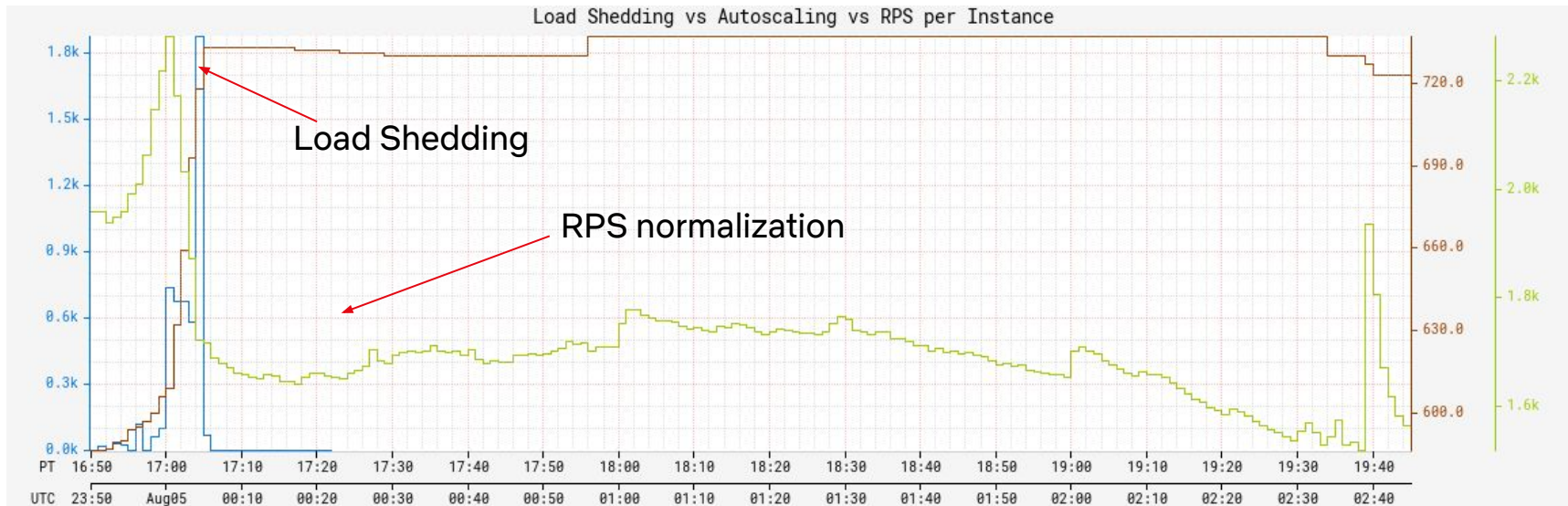
Max :	737.000	Min :	586.000
Avg :	724.246	Last :	723.000
Tot :	126.743k	Cnt :	175.000

Axis 2

RPS per Instance

Max :	2.280k	Min :	1.518k
Avg :	1.710k	Last :	1.565k
Tot :	299.236k	Cnt :	175.000

Fast Recovery Mitigates Pre-Scale Risk



Axis 0

Shed RPS

Max :	1.872k	Min :	0.000
Avg :	156.689	Last :	0.000
Tot :	5.014k	Cnt :	32.000

Axis 1

Fleet Size

Max :	737.000	Min :	586.000
Avg :	724.246	Last :	723.000
Tot :	126.743k	Cnt :	175.000

Axis 2

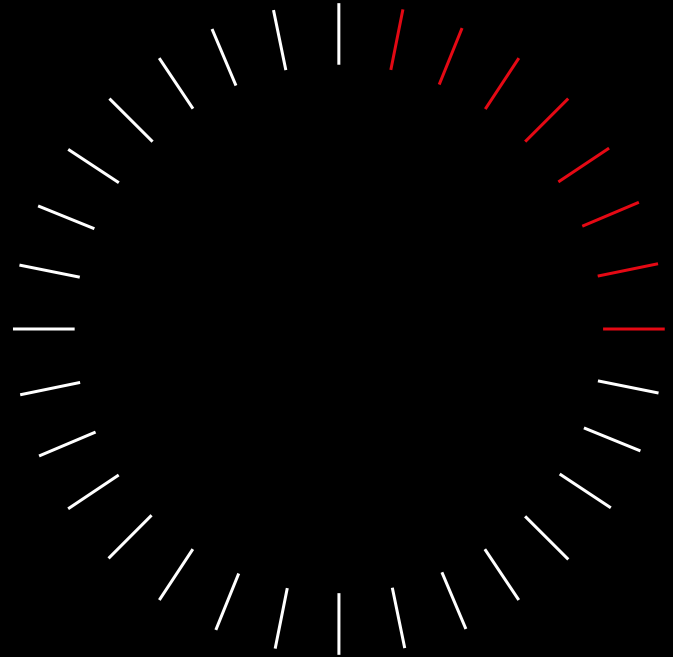
RPS per Instance

Max :	2.280k	Min :	1.518k
Avg :	1.710k	Last :	1.565k
Tot :	299.236k	Cnt :	175.000

Manage Demand

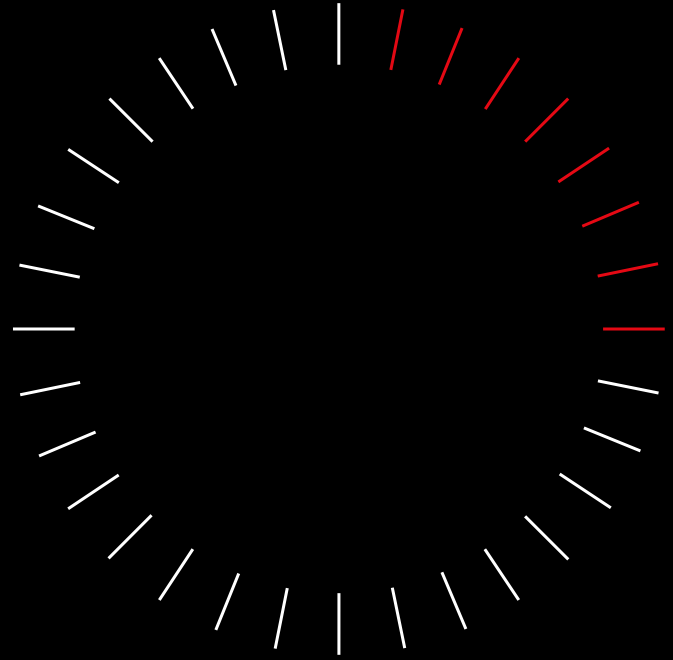
Pre-scaling

Dynamic Traffic Shaping

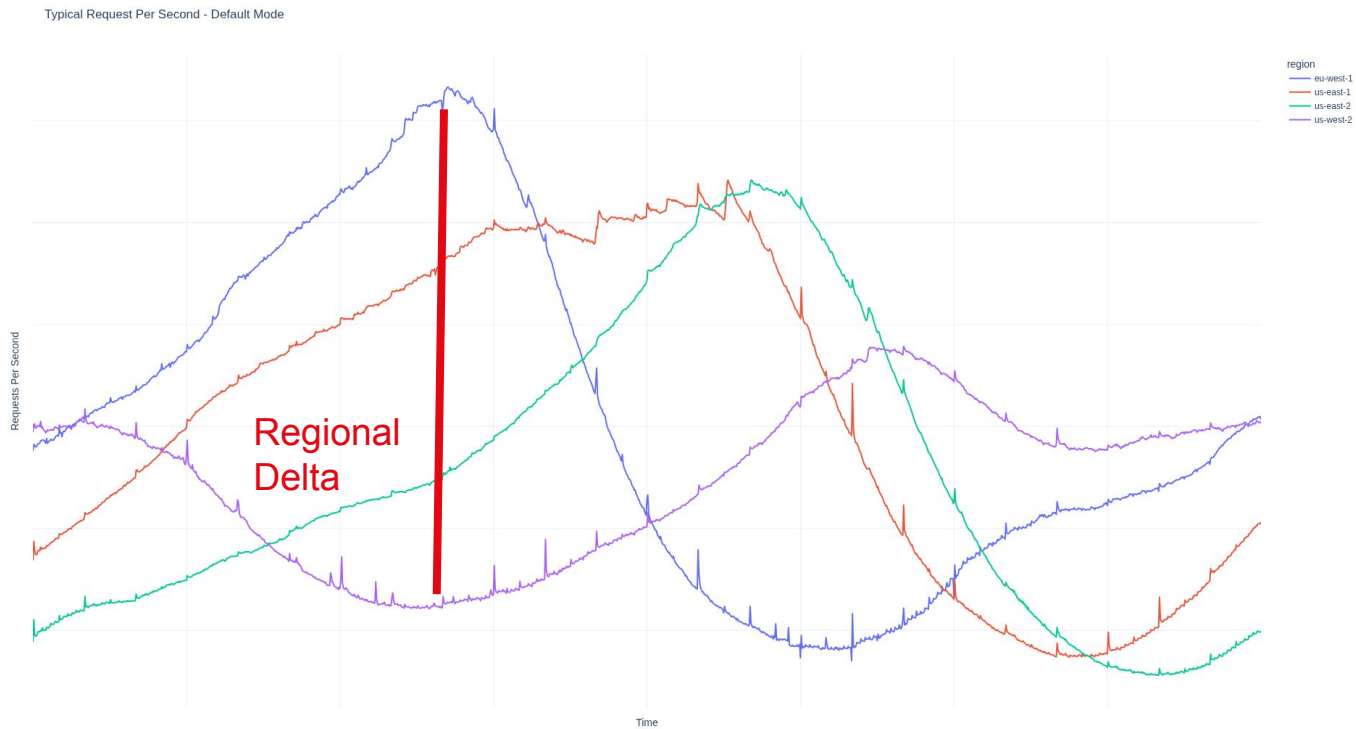


Dynamic Traffic Shaping

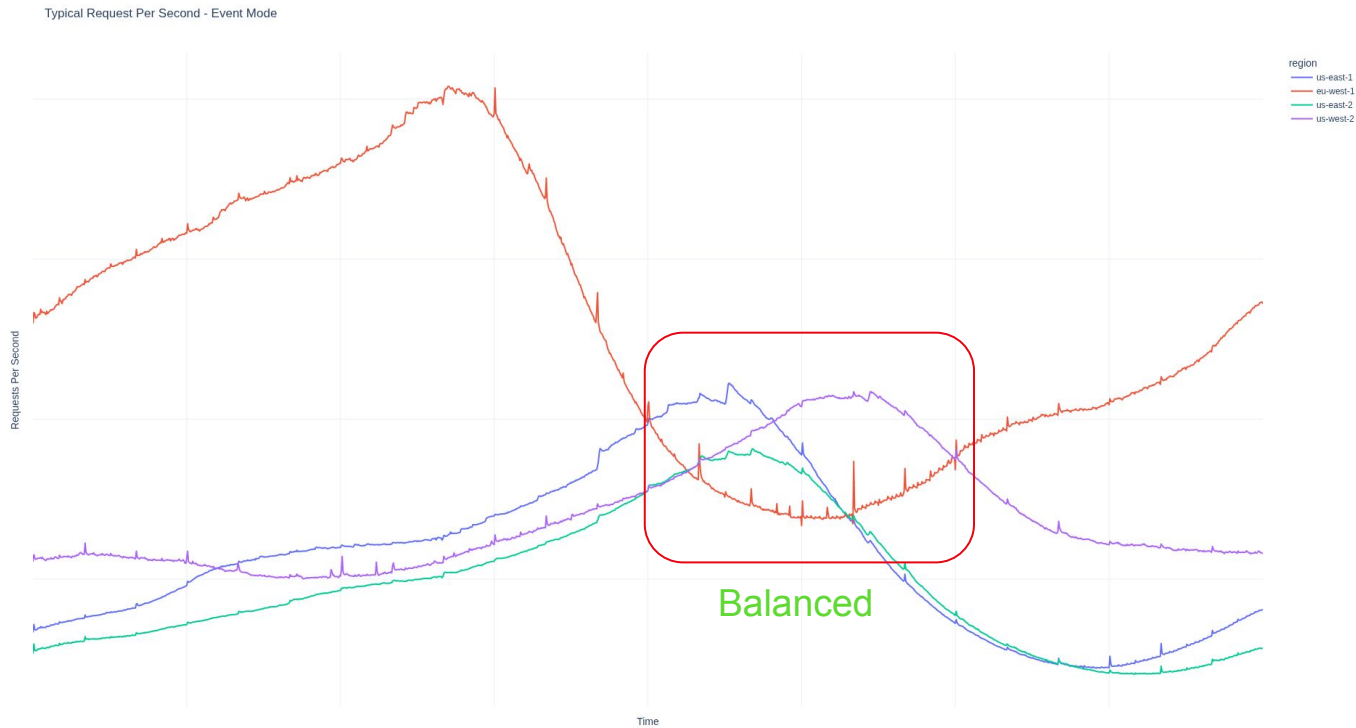
Re-distribute **existing** traffic



Default Shape



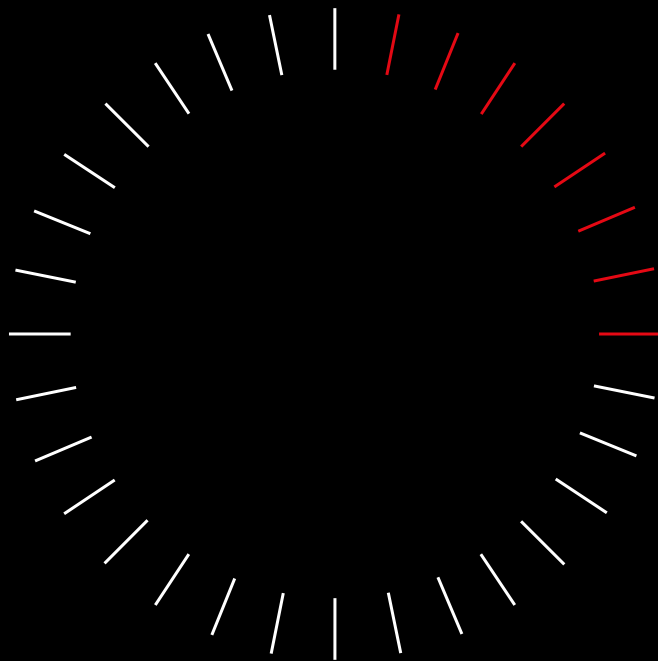
Shaping in action



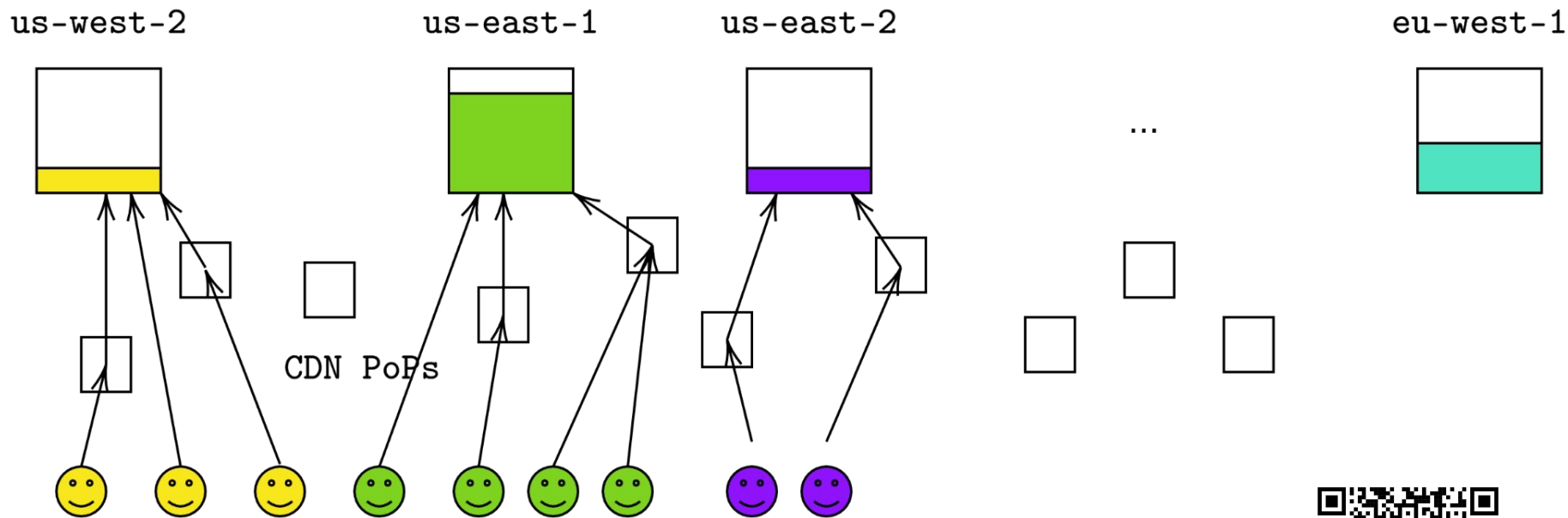
Dynamic Traffic Shaping

Re-distribute **existing** traffic 

Steer **new** traffic based on latency



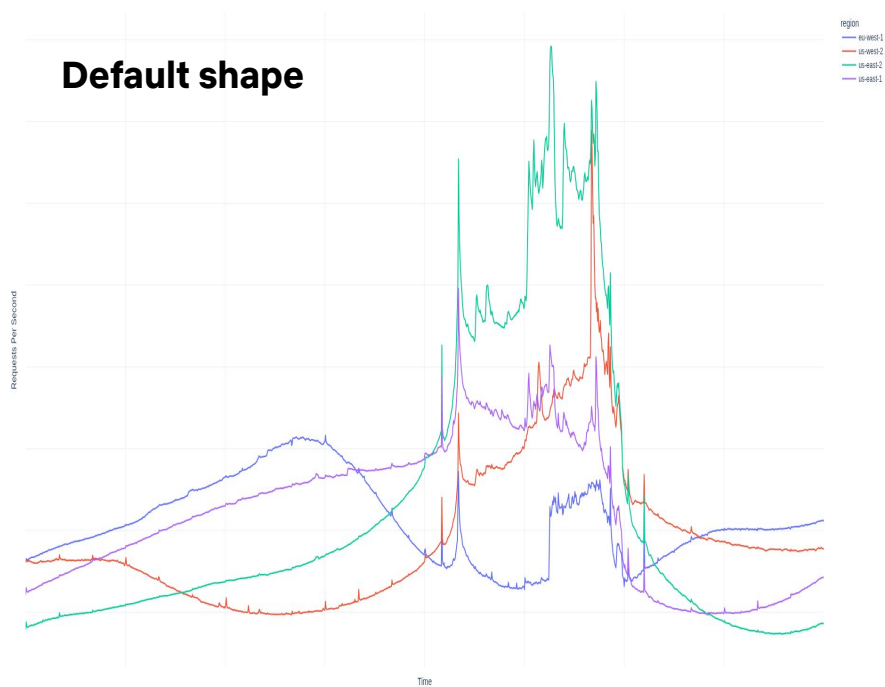
Netflix DNS and Steering - Typical



Shaping effects

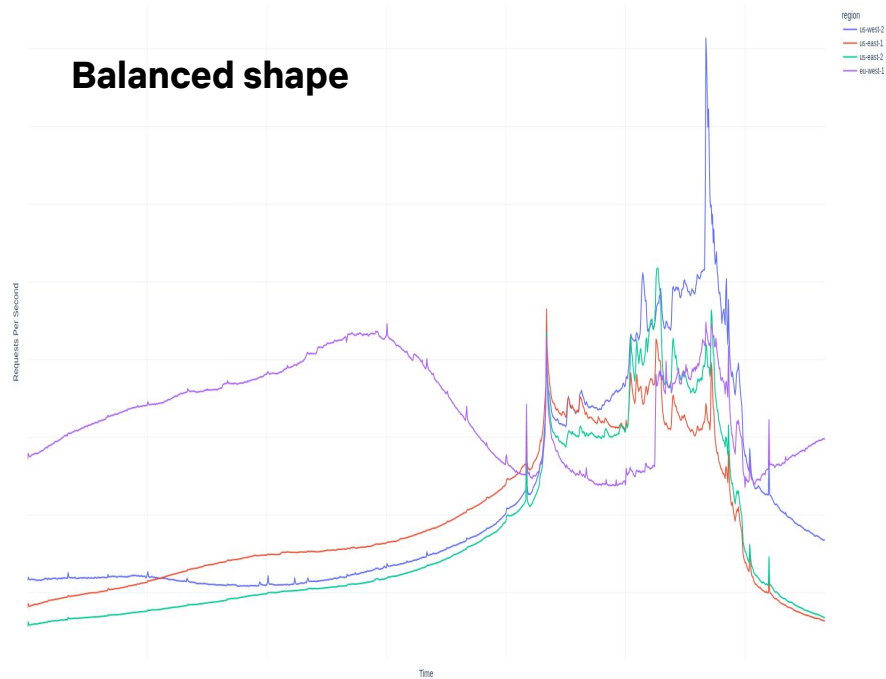
Event Request Per Second - Default Mode

Default shape



Event Request Per Second - Event Mode

Balanced shape



Hope is not a *Strategy*...



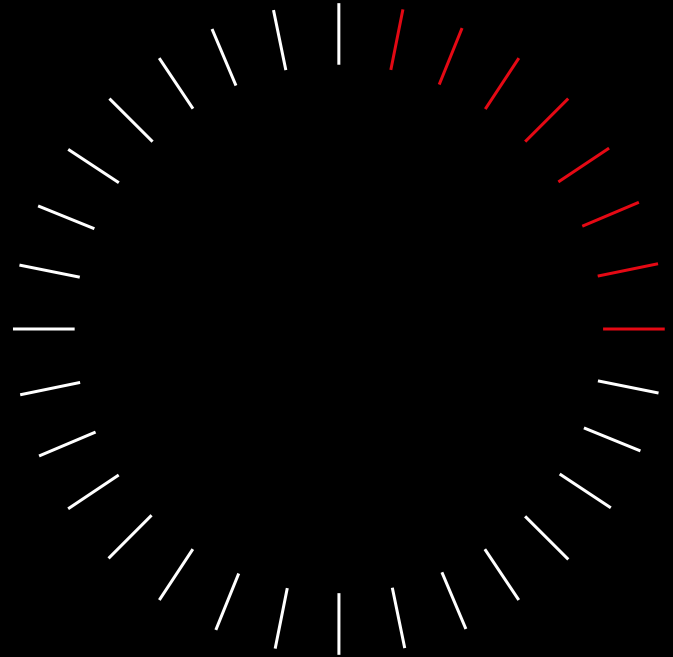
Perfectly Match
Supply with Demand

Manage Demand

Pre-scaling ✓

Dynamic Traffic Shaping ✓

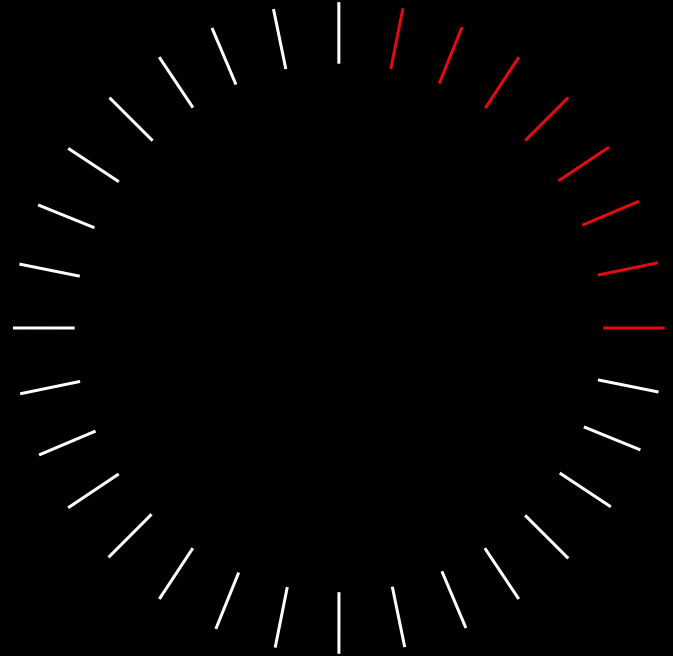
Reactive Scaling



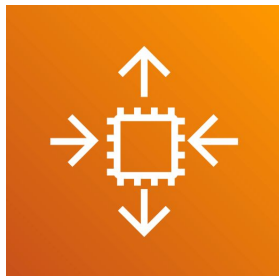
Reactive Scaling

Buffers based on workload tier/domain

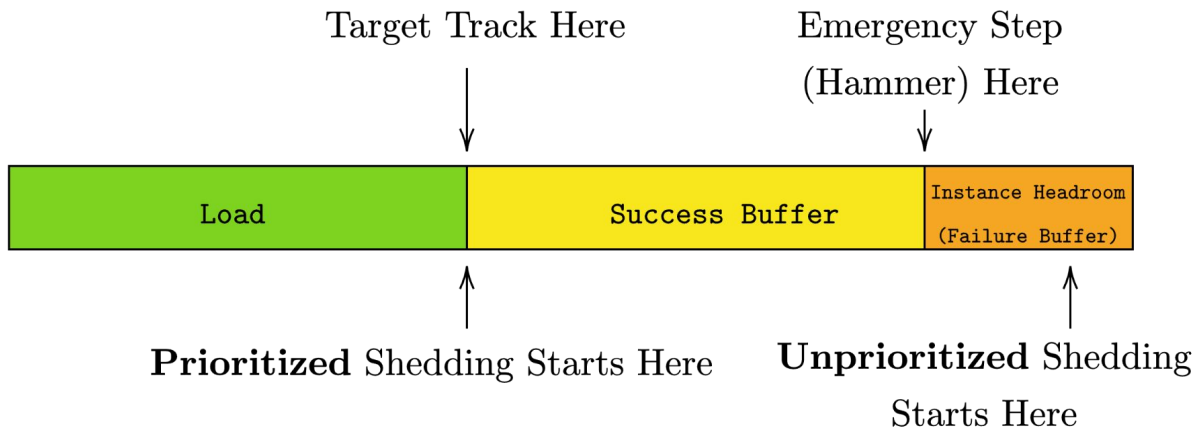
Mathematically derived **thresholds**



Scaling Policies

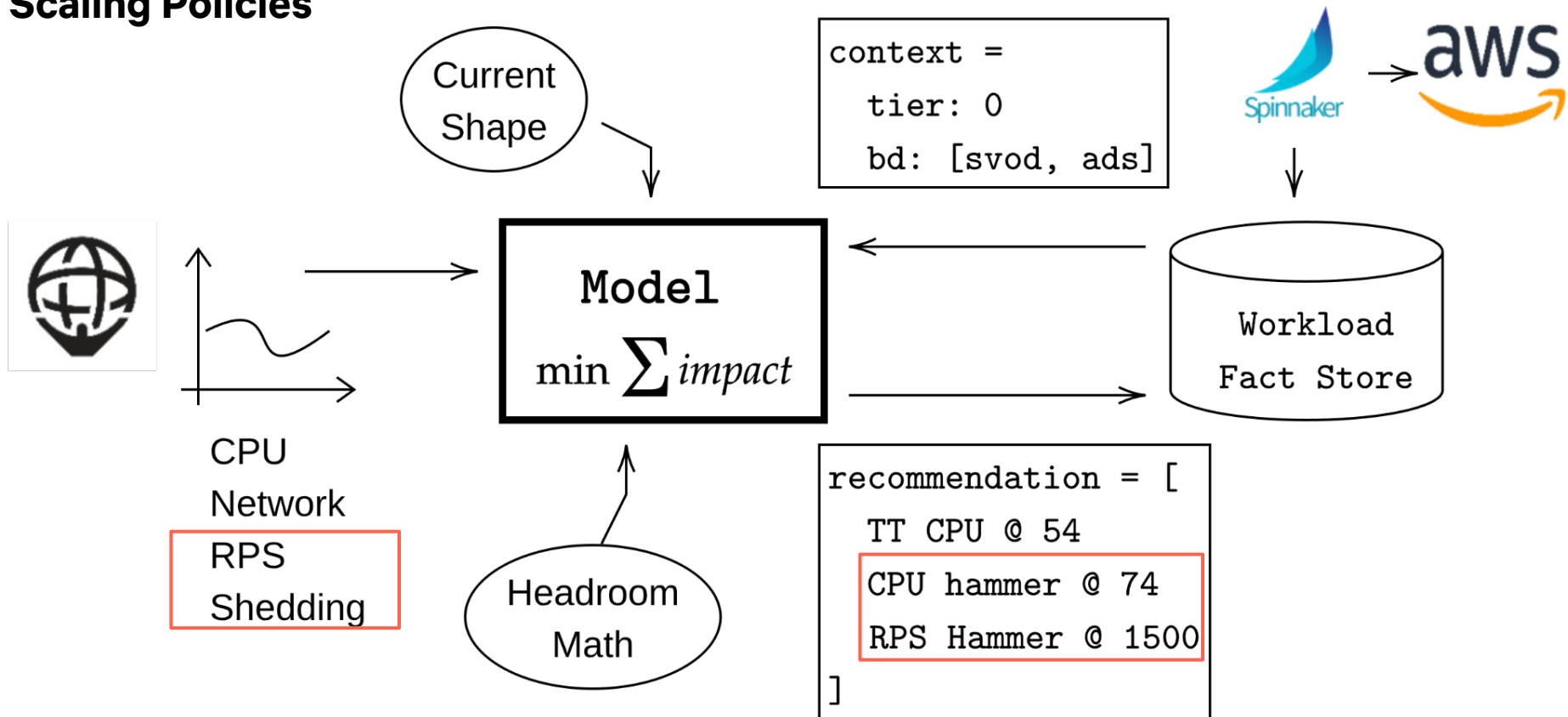


CPU/Network **Autoscaling** Policy



Load Shedding Policy

Scaling Policies

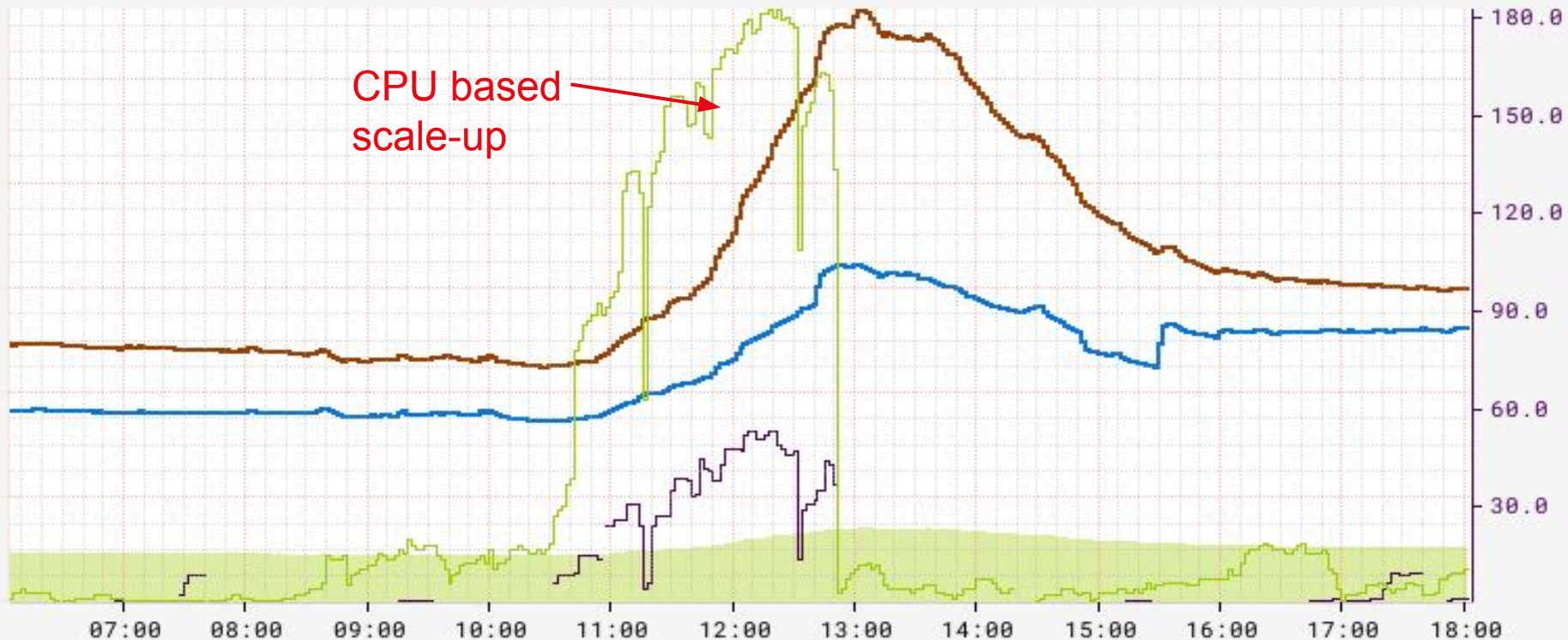


Reactive Scaling in action

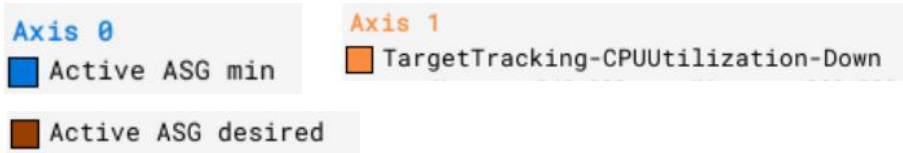
Active ASG desired
Active ASG min

TargetTracking-CPUUtilization-Up

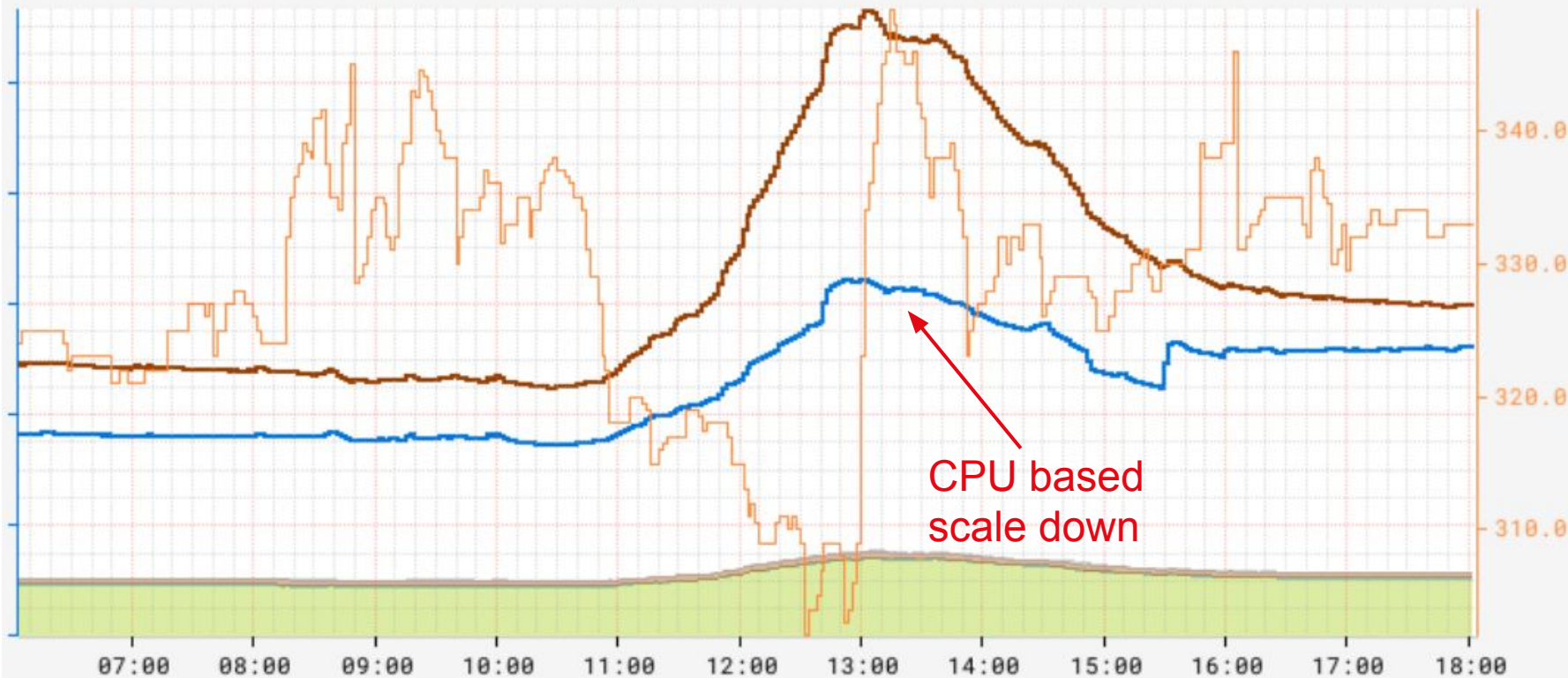
Cluster Footprint - Scale Down Events



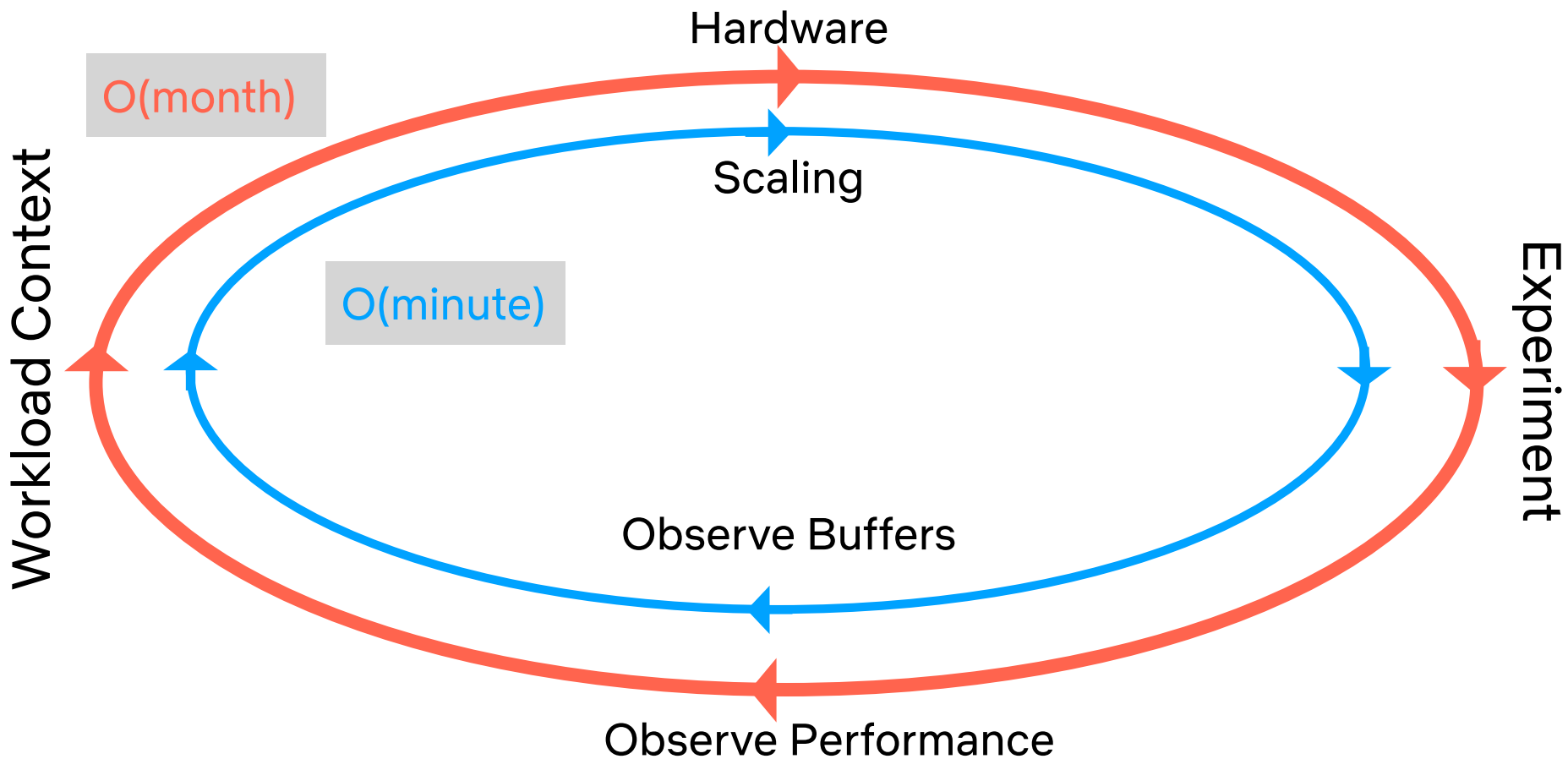
Reactive Scaling in action



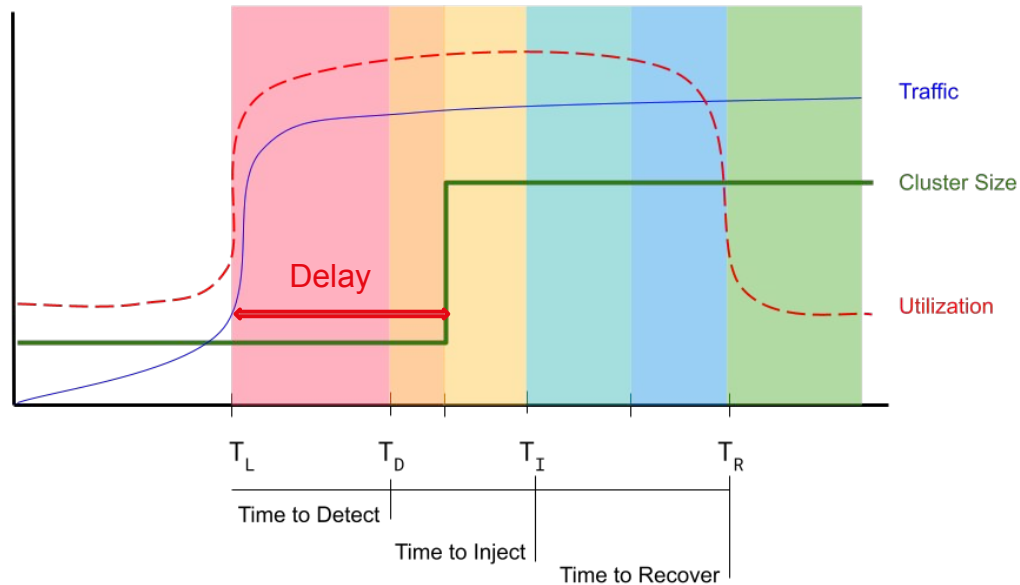
Cluster Footprint - Scale Down Events



CPU based
scale down



But....auto-scaling 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 .



Does reactive *FAST* scaleup work in practice?



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

8-15 mins

Does reactive *FAST* scaleup work in practice?



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

8-15 mins

~70% reduction



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

3-4 mins

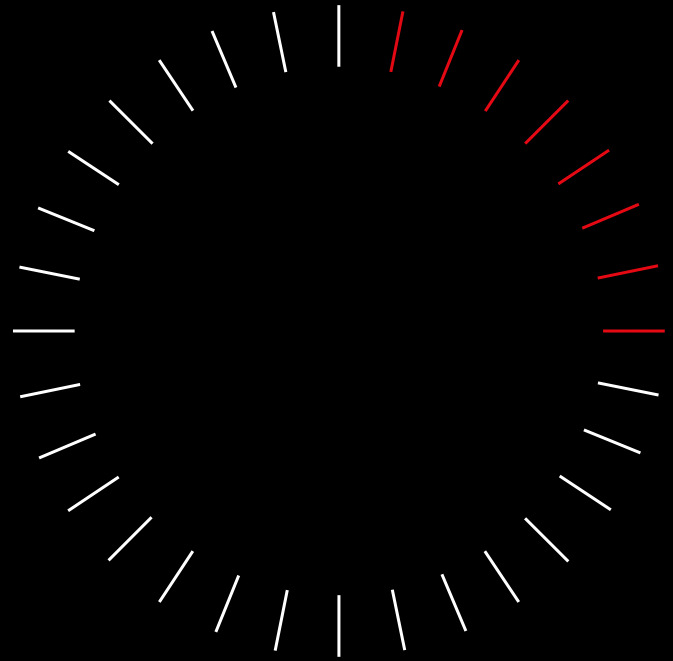
Manage Demand

Pre-scaling

Dynamic Traffic Shaping

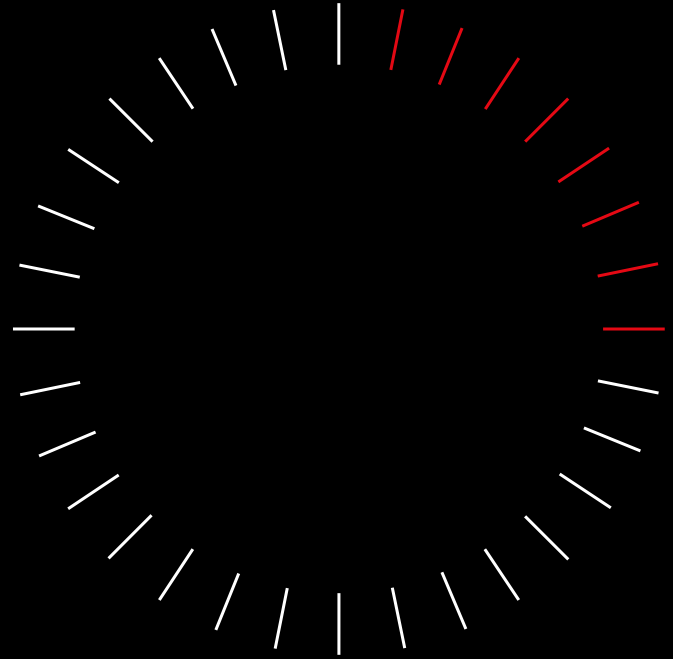
Reactive Auto-Scaling

Load Shedding



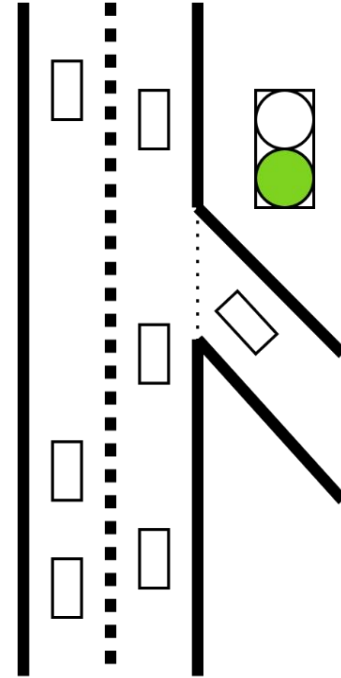
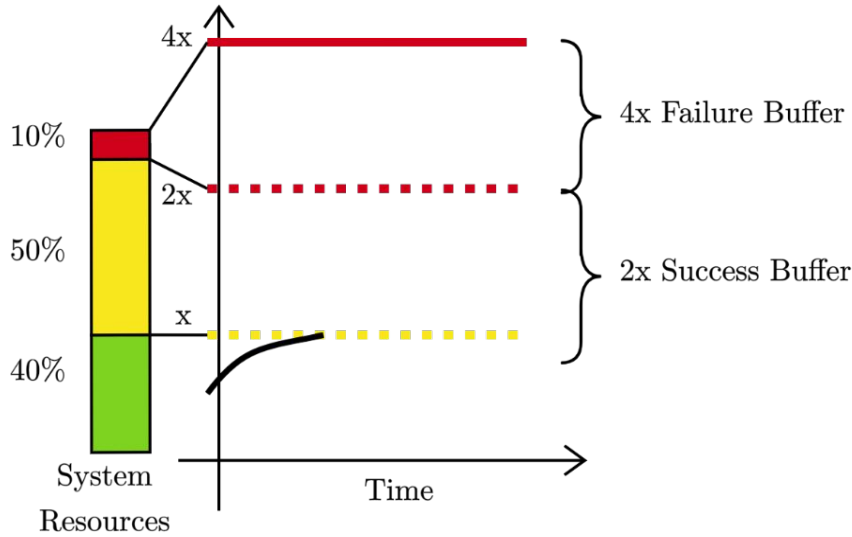
Load Shedding

Bulk load shedding



CPU Load Shedding

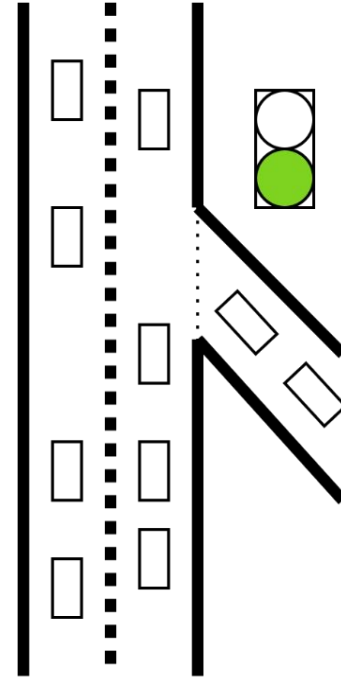
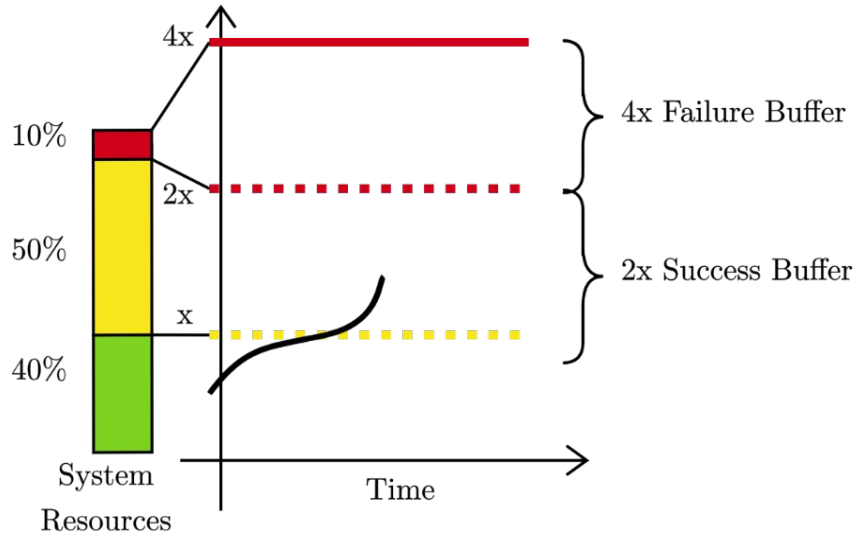
Normal System Load with Buffer



$$T_{transit} = 10m$$

CPU Load Shedding

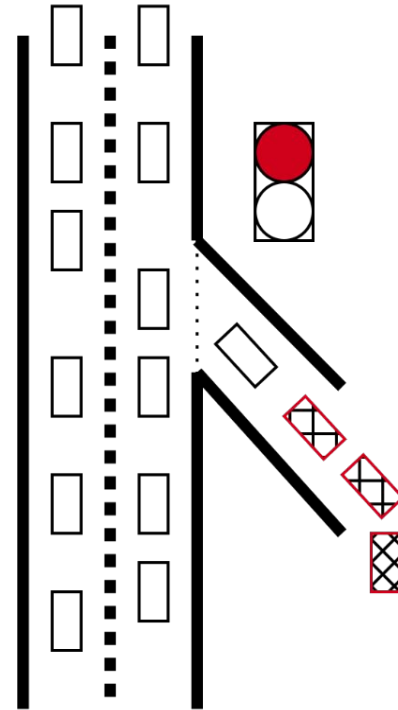
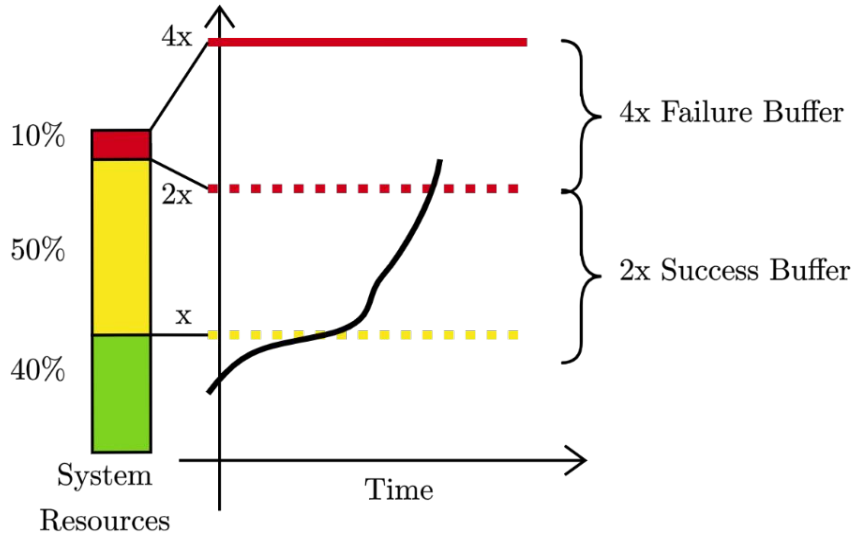
Normal System Load with Buffer



$$T_{transit} = 15m$$

CPU Load Shedding

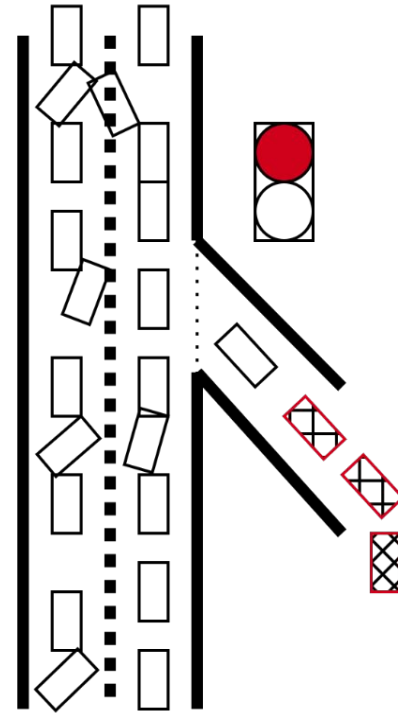
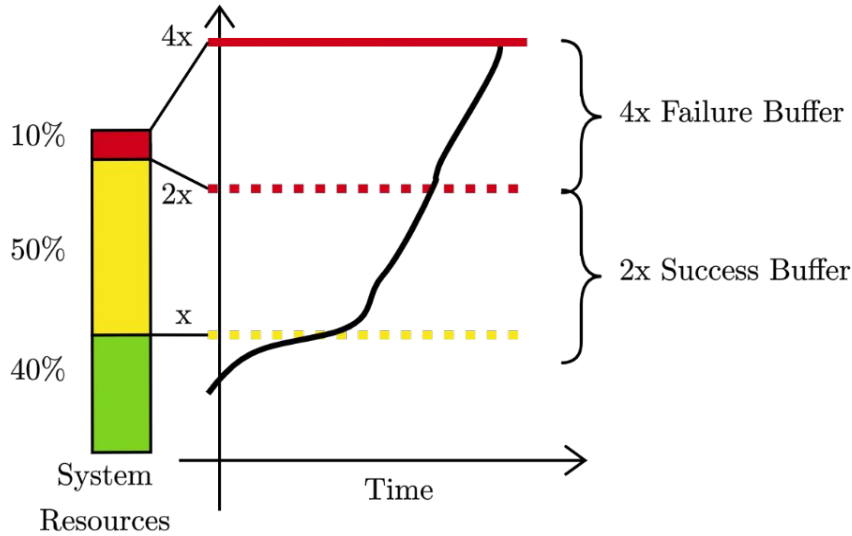
Normal System Load with Buffer



$$T_{transit} = 30m$$

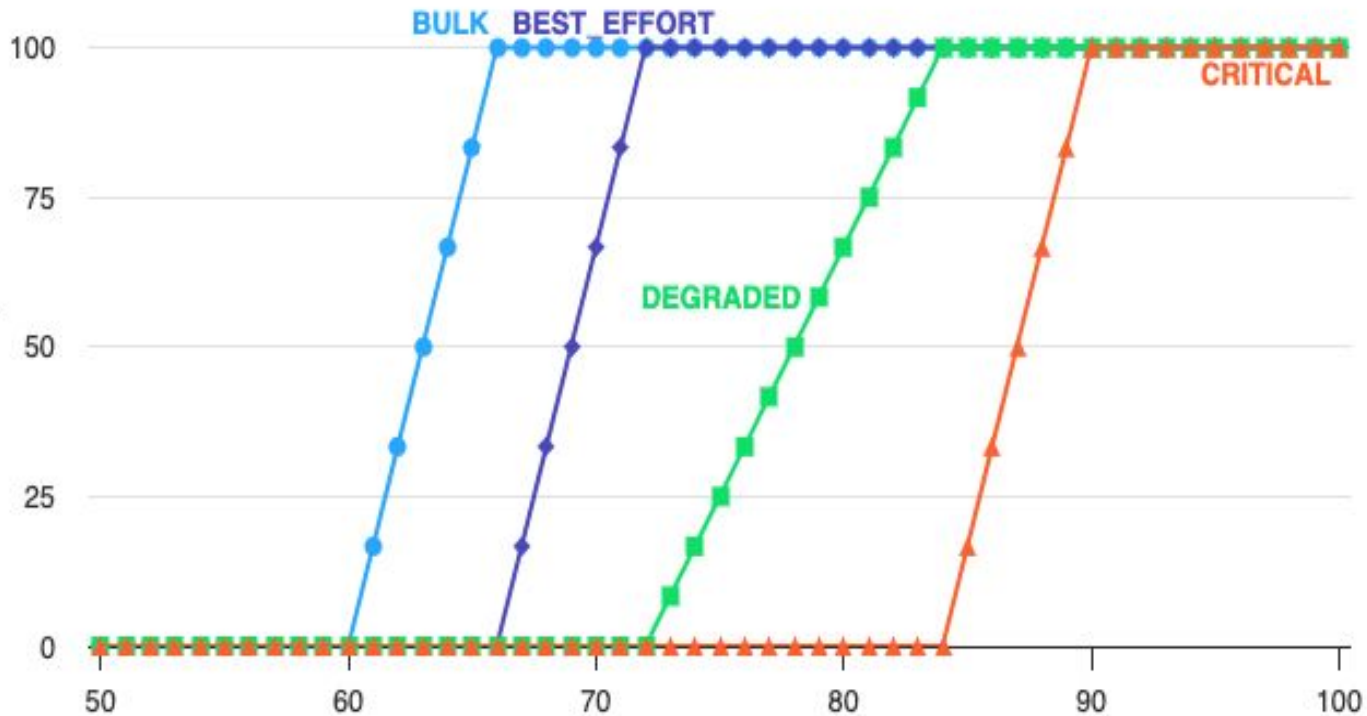
CPU Load Shedding

Normal System Load with Buffer



$$T_{transit} = \infty$$

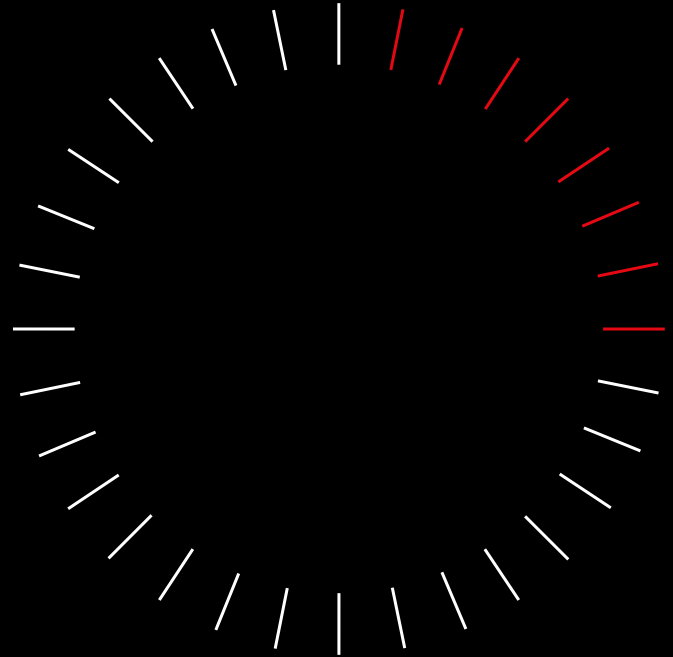
CPU Load Shedding



Load Shedding

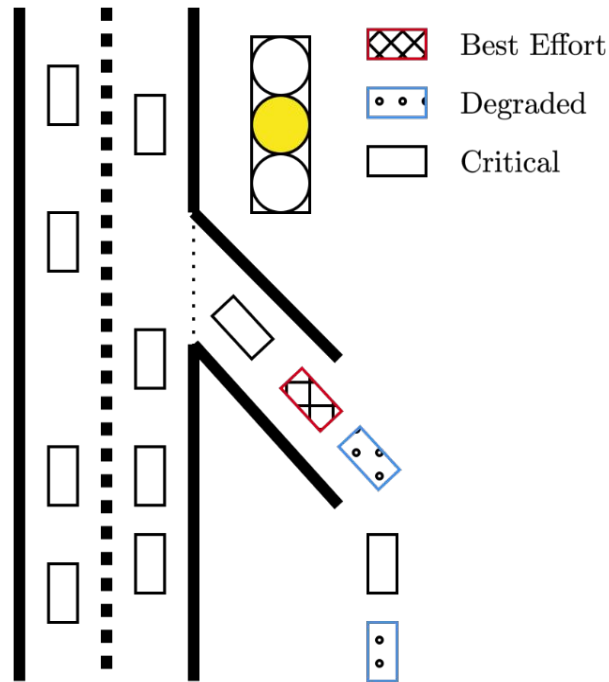
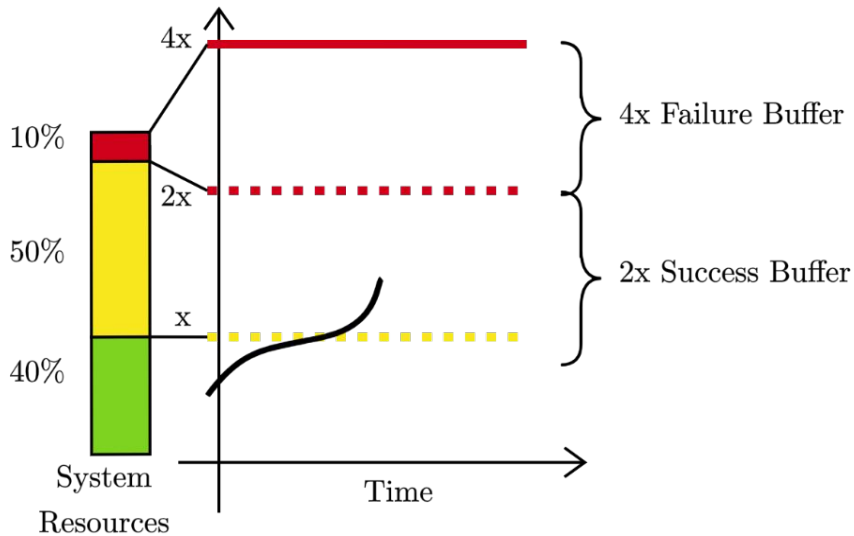
Bulk load shedding ✓

Priority based load shedding



Prioritization creates additional Buffer

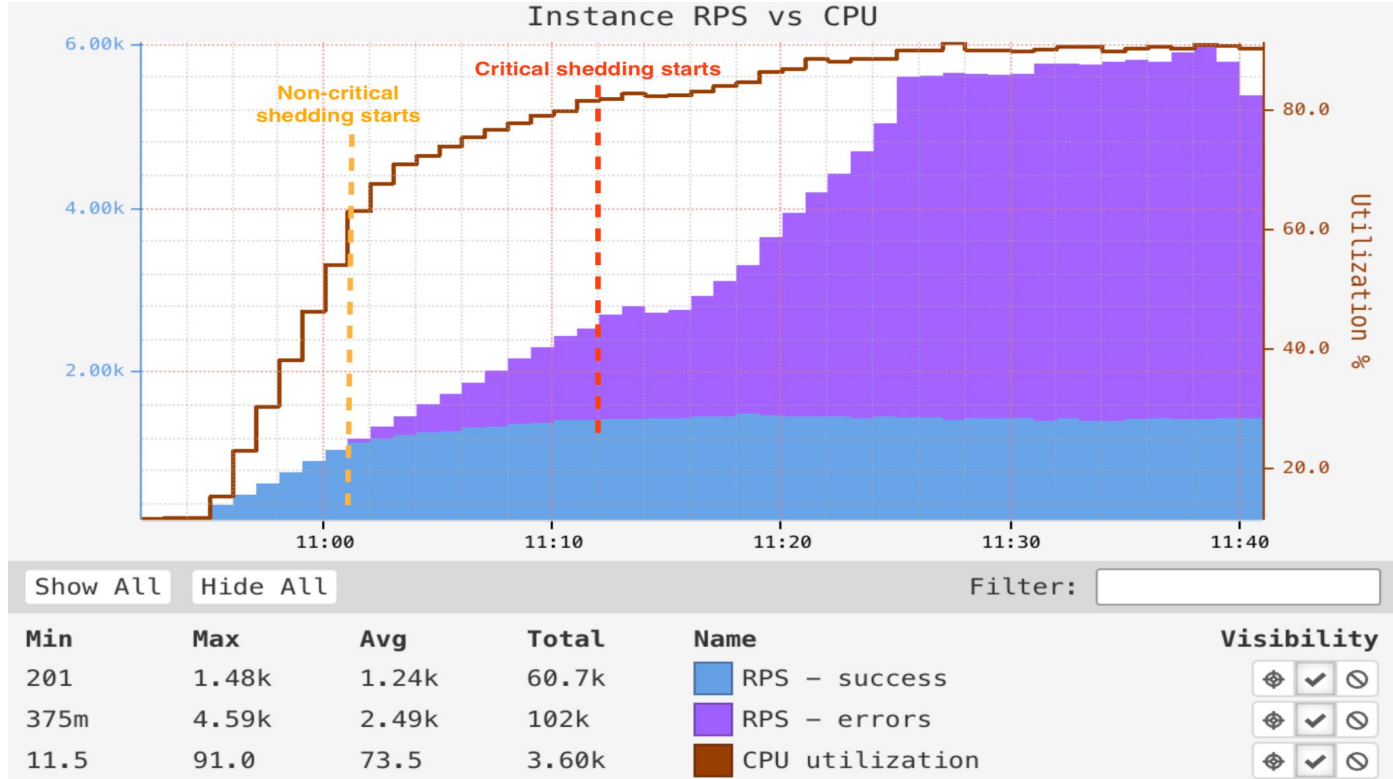
Normal System Load with Buffer

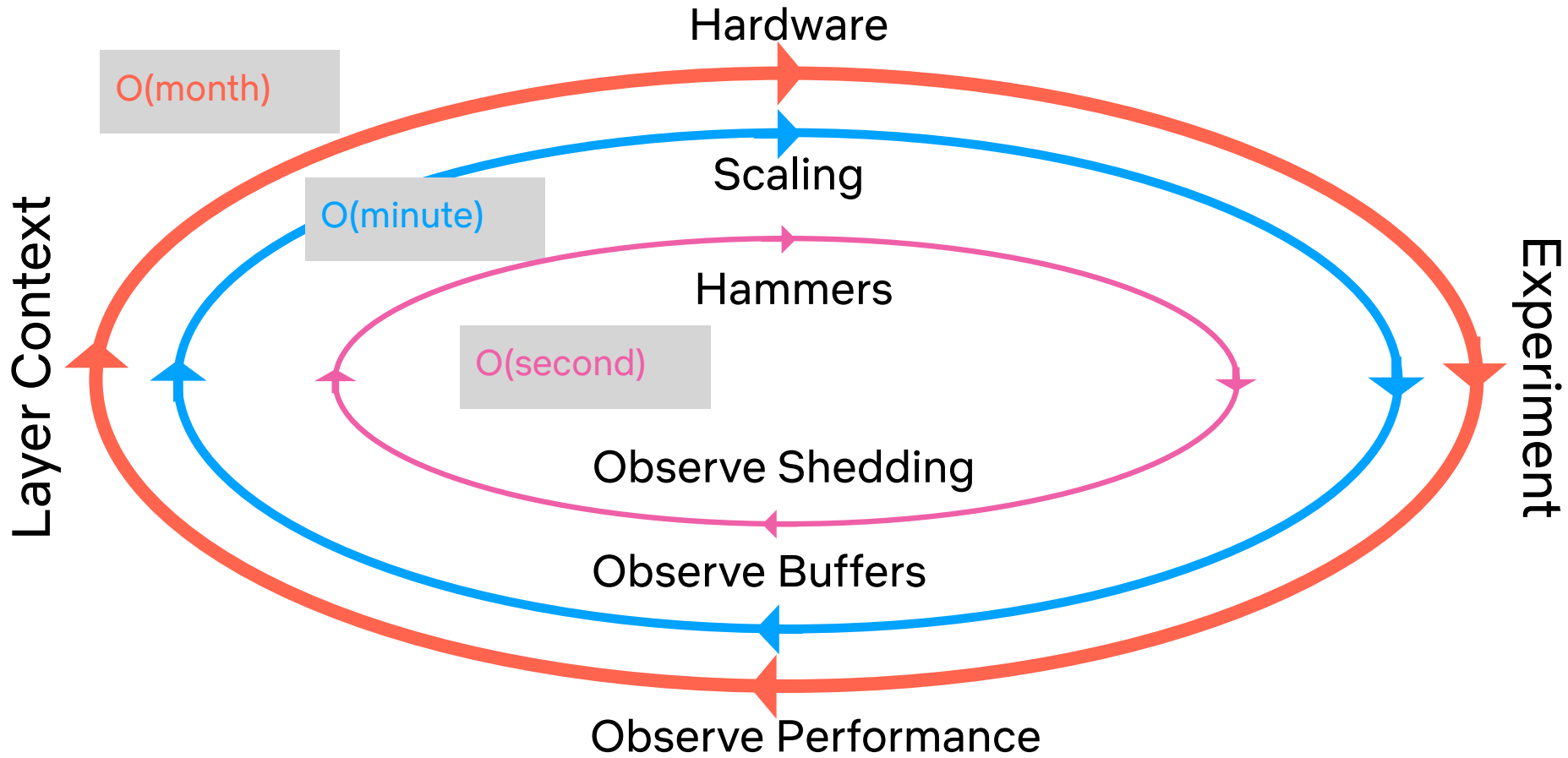


$$T_{transit} = 15m$$



Load Shedding Priority based





Lesson #1

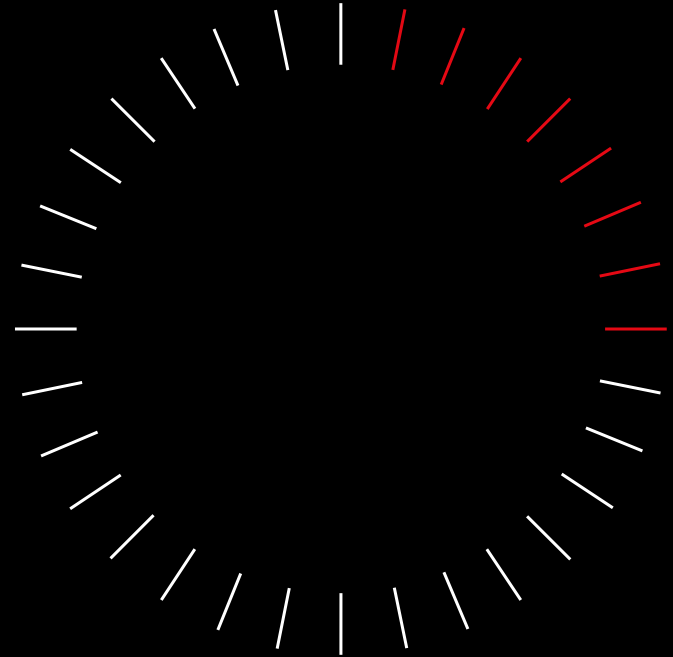
Proactive **and** Reactive levers

Pre-Ingress

- Predictive Scaling
- Traffic Shaping

Post-Ingress

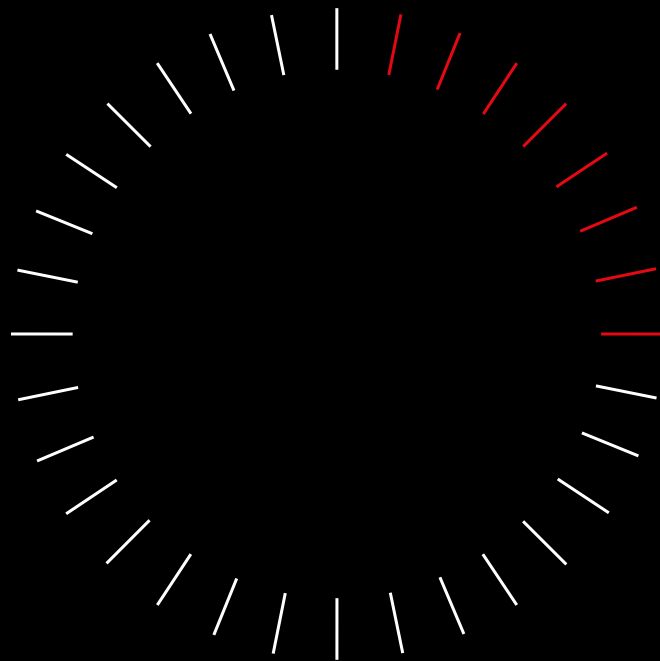
- Reactive Scaling
- Load Shedding



Lesson #2

Systems Thinking..

- Raw Compute isn't the silver bullet to buy down **Risk**
- Efficiency wins **compound**, at scale.
- **End to End** Traffic Management

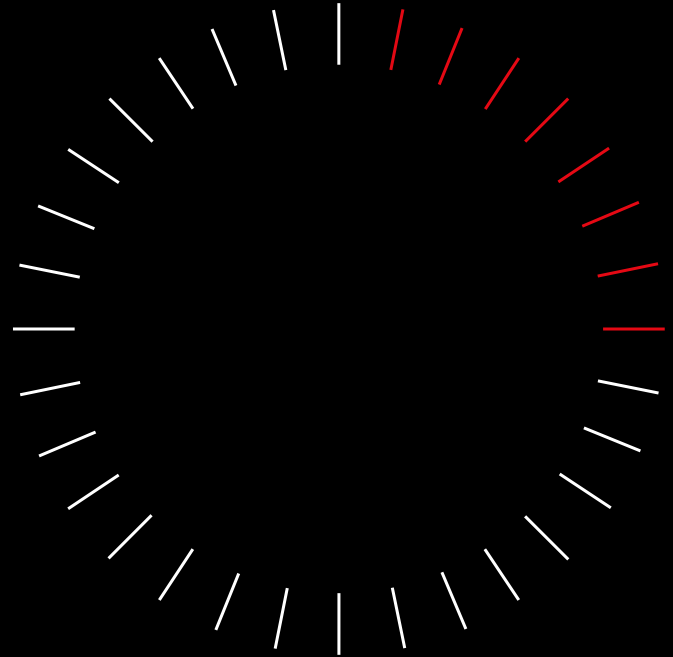


Lesson #3

Math is our Safety Blanket™



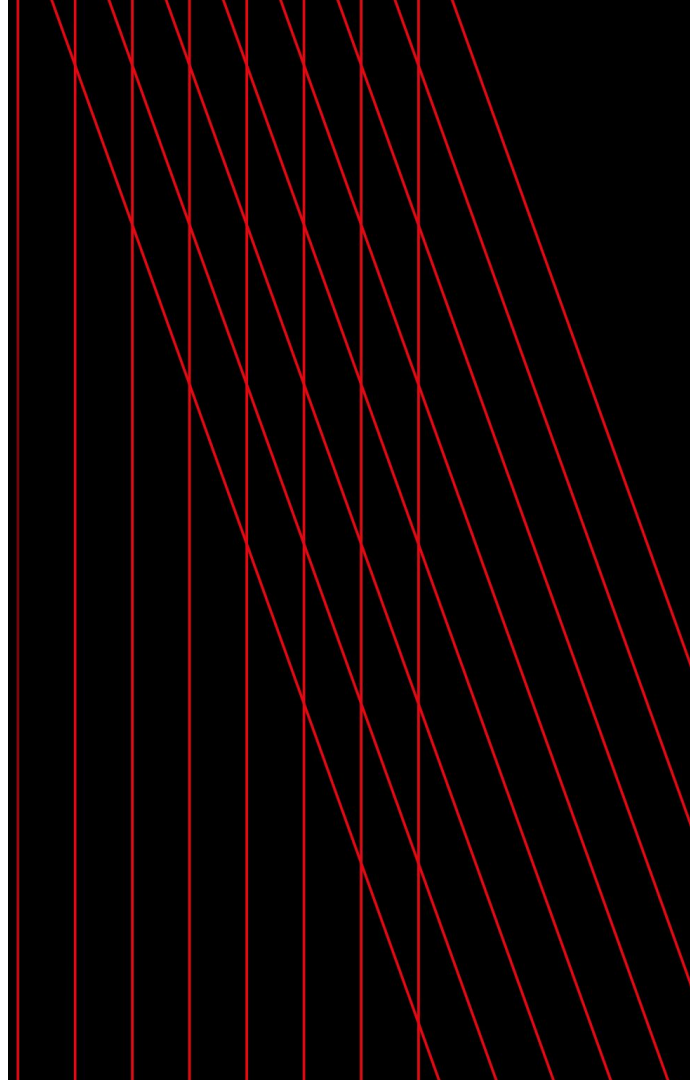
Econ 101



Thank You



Argha C, arghac@netflix.com
Joseph Lynch, josephl@netflix.com



Please rate and leave feedback!



Argha C, arghac@netflix.com
Joseph Lynch, josephl@netflix.com

Questions?