# Capacity vs Efficiency

Building a Globally Scalable Cloud Database

SRECon 2022, Sydney          Daniel Marshall <djmarshall@google.com>

# Me

Electrical engineer turned SRE at
Google Sydney

I work on Google Cloud Firestore

If you want to chat more about this
topic, come find me after

# Capacity **vs** Efficiency

**Capacity** – provisioning sufficient resources to handle incoming request load

**Efficiency** – minimising the total cost of the system

Google Cloud Firestore is a **fully scalable** NoSQL database where you **only pay for what you use**.

Corollaries

- Fully scalable ⇒ we must be able to serve any amount of traffic, at any time

- Pay for use ⇒ we pay for everything else

# Capacity

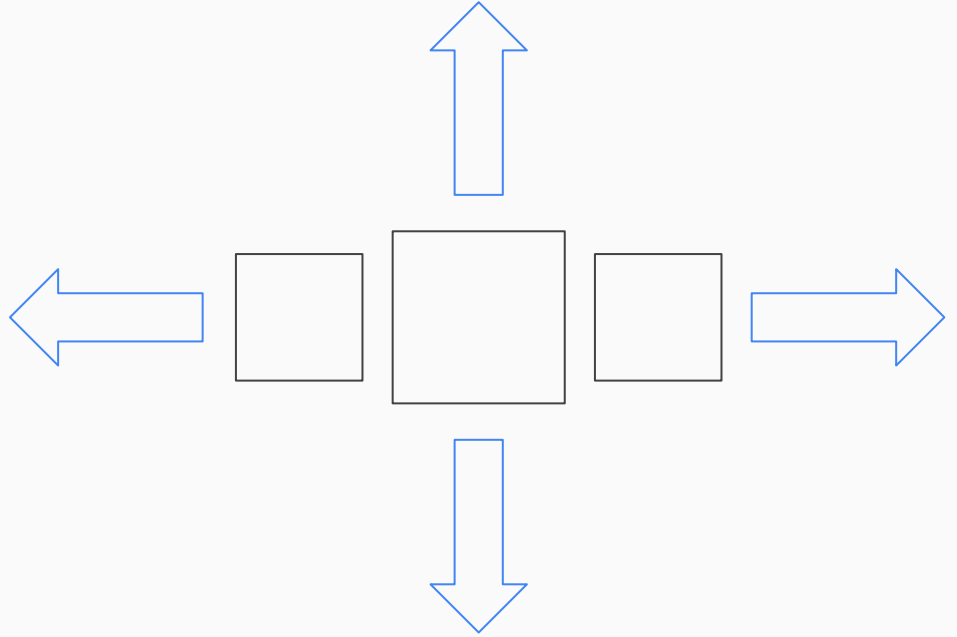Autoscaling is the first step

Improves reliability **and** toil
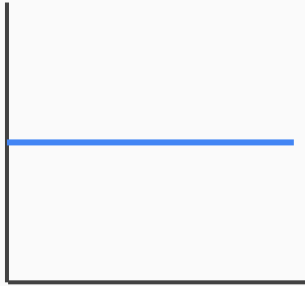
Autoscale horizontally and vertically

Choose one metric to scale on
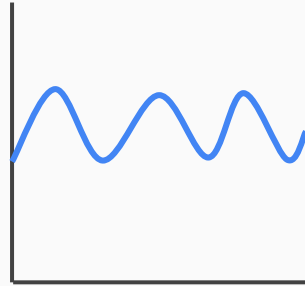
Align your bottlenecks with your scaling signal

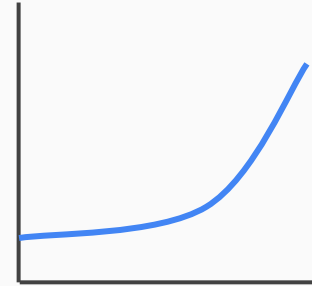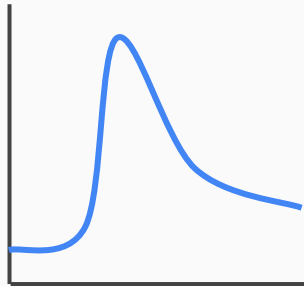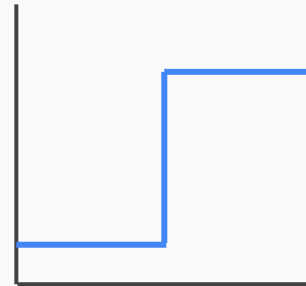(Marginally) overprovision everything else

# Traffic Patterns

### Flat

### Diurnal

### Slow ramp-up

### Traffic spike

### Bulk workload

Where possible, **control the ramp-up** speed

# Choosing Your Headroom

Reaction Window

Headroom

Capacity

Load

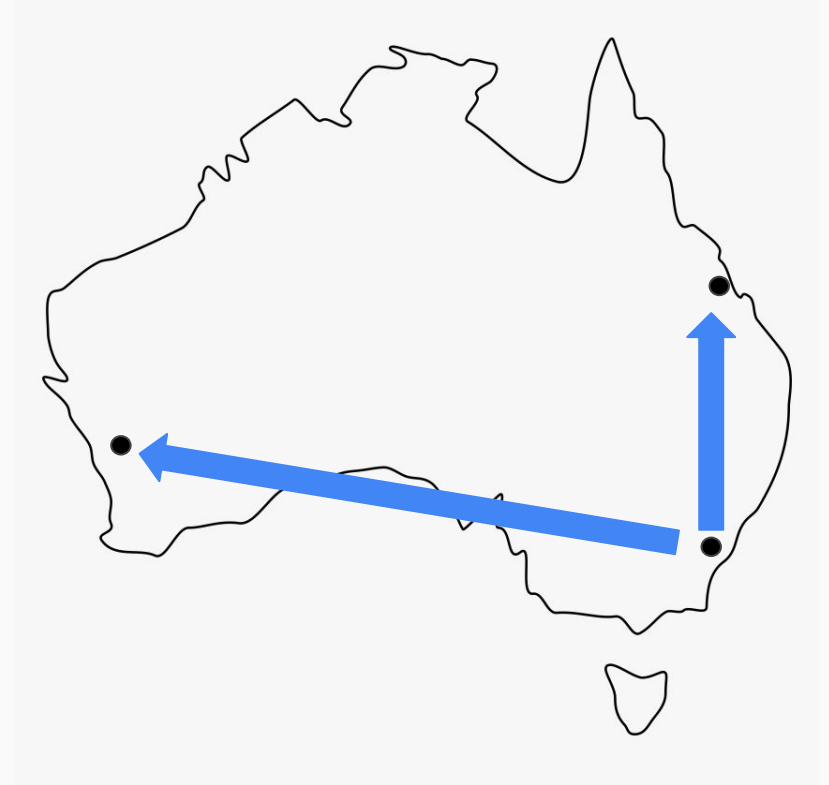Efficiency = 100% − max(traffic_spike)

# Stockout Resilience

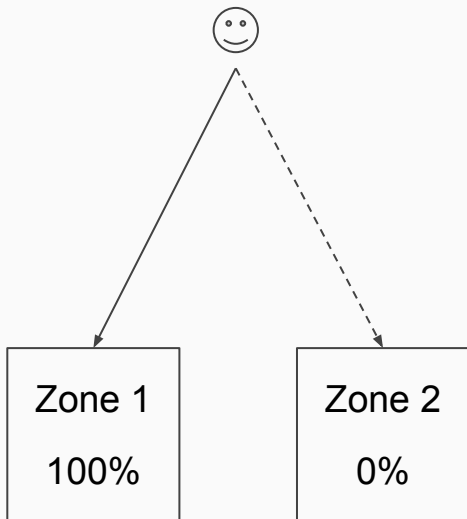Stockout: underlying platform runs out of capacity

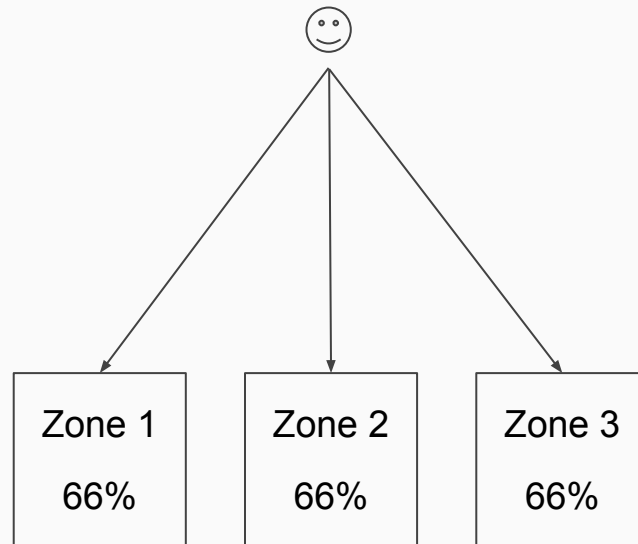Use as many availability zones and regions as you can

Build capacity agility
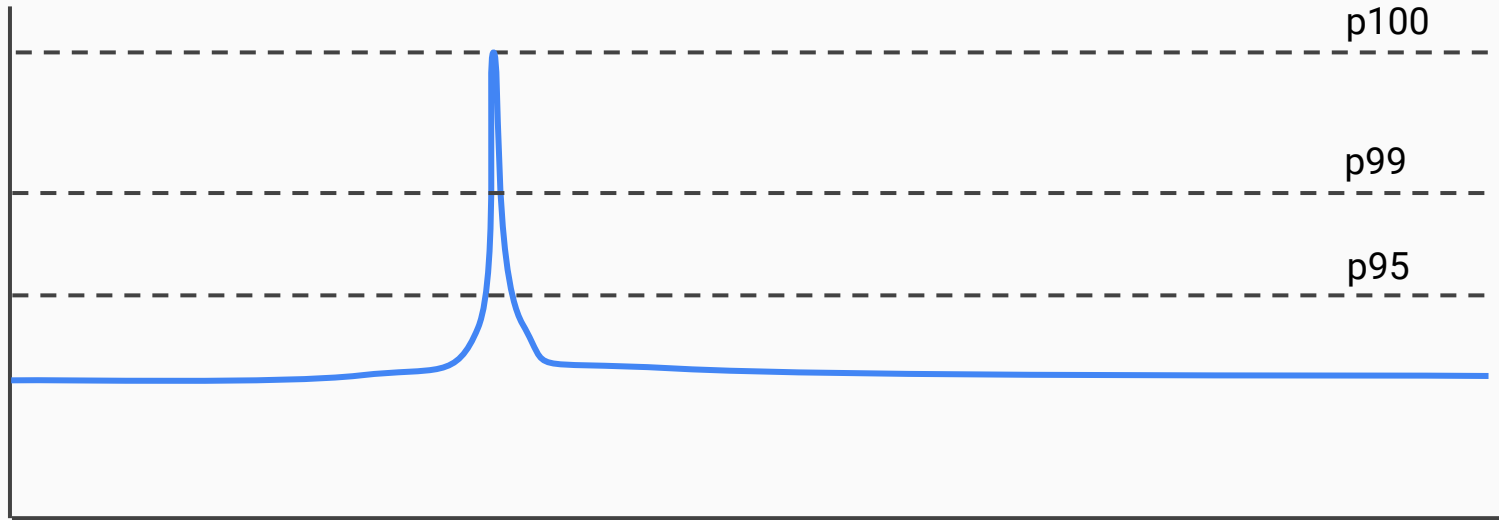
# Provisioning For Failover

Active/passive

```
        ☺
       / \
      /   \
     ↓     ↓
┌────────┐ ┌────────┐
│ Zone 1 │ │ Zone 2 │
│        │ │        │
│  100%  │ │   0%   │
└────────┘ └────────┘
```

N+1

```
         ☺
       / | \
      /  |  \
     ↓   ↓   ↓
┌──────┐┌──────┐┌──────┐
│Zone 1││Zone 2││Zone 3│
│      ││      ││      │
│ 66%  ││ 66%  ││ 66%  │
└──────┘└──────┘└──────┘
```

All failover models are wrong. Test your failovers!

# Efficiency

## Autoscaling Target

# Tuning The Autoscaling



Reaction Window

Stabilization / decay period

Capacity

Load

# Tuning The Autoscaling

# Load Balancing



Good load balancing

Poor load balancing

Max

Mean

Load

Servers

Efficiency = 100% − max(traffic_spike) **− load_imbalance**$_{max-mean}$

Efficiency = 100% − max(traffic_spike) + load_imbalance$_{\text{max-mean}}$



Efficiency = 100% − $\mathbf{P_{max}}$(traffic_spike + load_imbalance$_{\text{max-mean}}$)
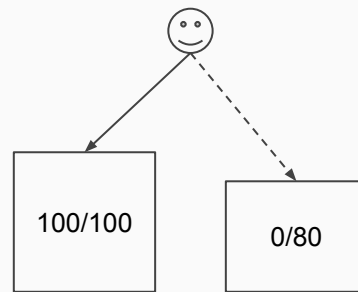
N+1: Increase the N!

- N=1    50% efficiency loss
- N=4    20% efficiency loss

Active/Passive

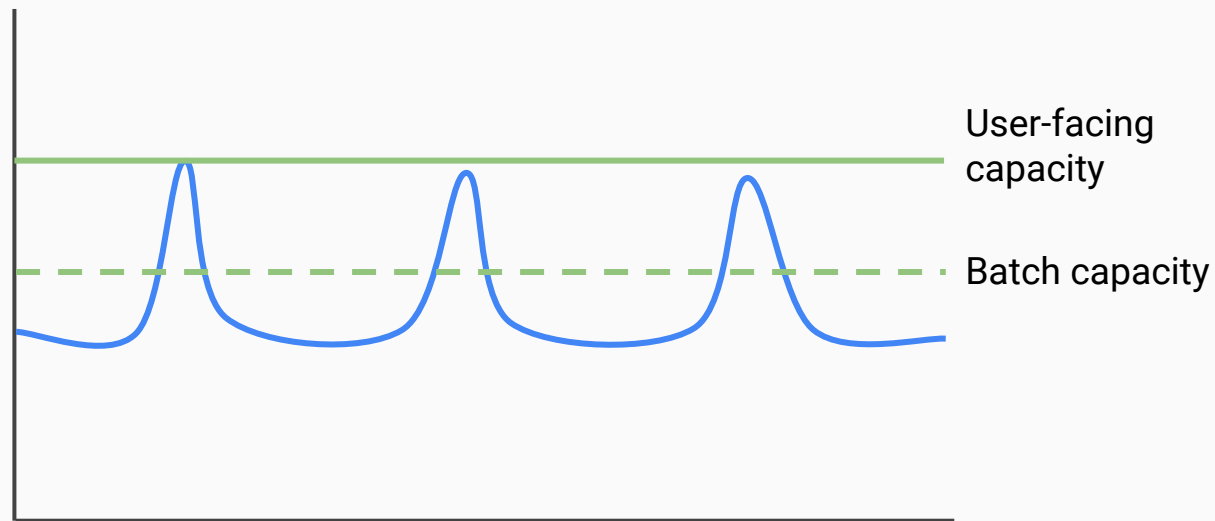- Use lower SLO / spot instances
- Underprovision and scale

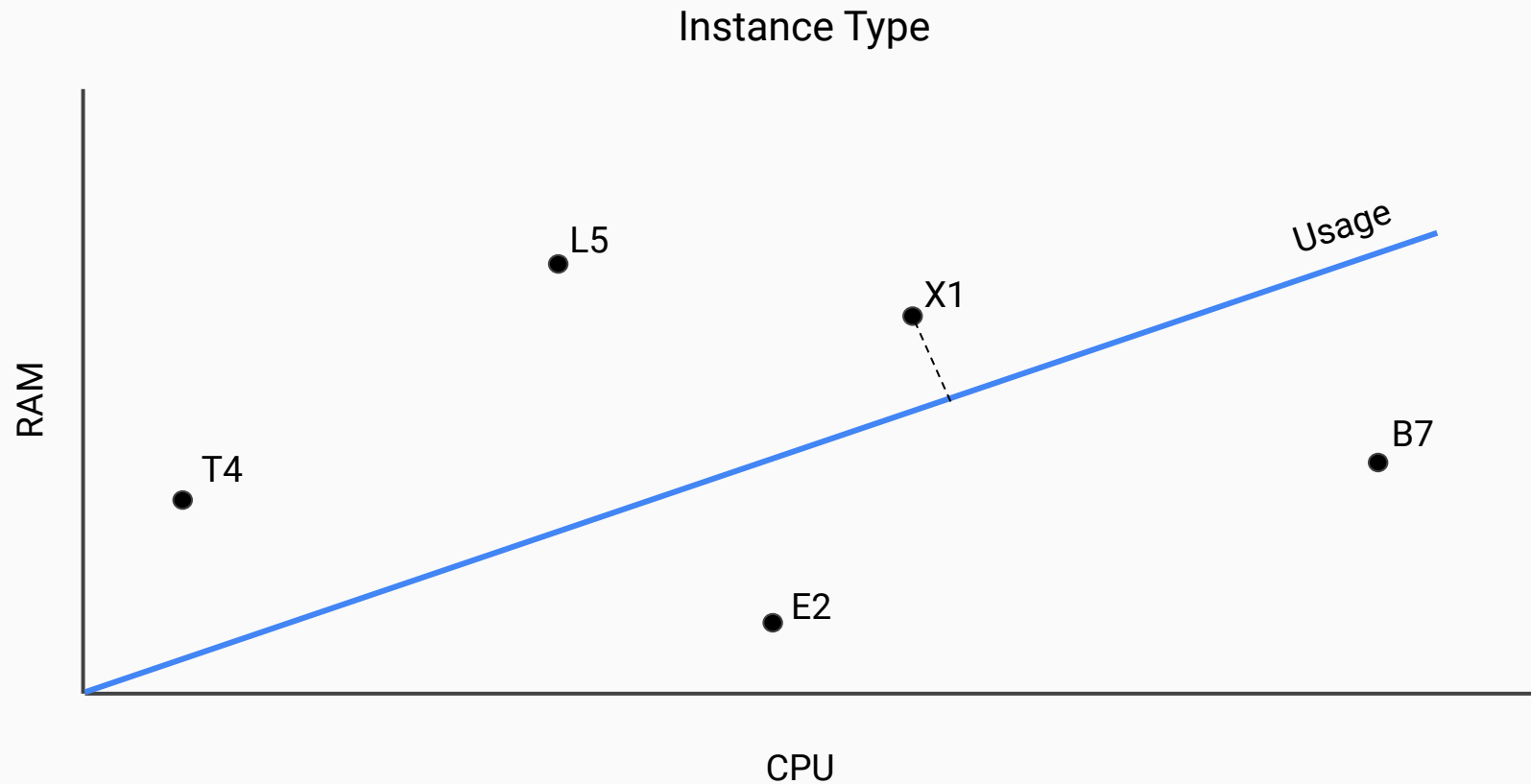Efficiency = 100% − $P_{max}$(traffic_spike + load_imbalance$_{max-mean}$ + **failover_capacity**)

# Batch Traffic

Key observation: Batch traffic is latency tolerant

- Split it out

- Use lower SLO
  / spot instances

- Run it hotter

User-facing capacity

Batch capacity
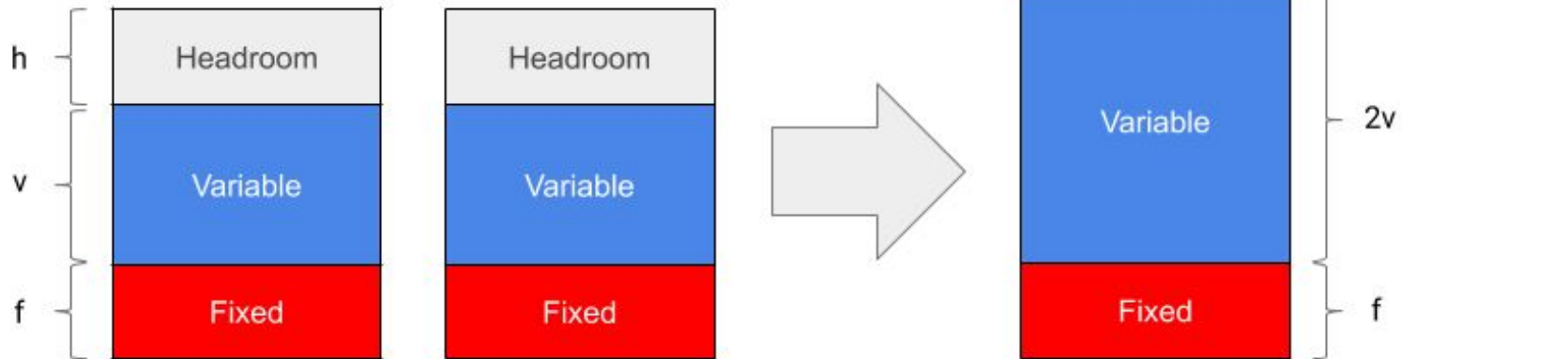
# Finding The Right Shape

# Bigger Is Better



Job with fewer, bigger tasks

h | Headroom | Variable | f | Fixed

f = Savings

<2h

2v

f

Byte/$:   HDD > SSD > RAM

IO/$:      HDD < SSD < RAM

For GC languages: CPU vs RAM



Bytes

IOPS

# Handling Overload

# Overload And Efficiency

Efficiency = **safe_cpu**% − $P_{max}$(traffic_spike + load_imbalance$_{max-mean}$ + failover_capacity)
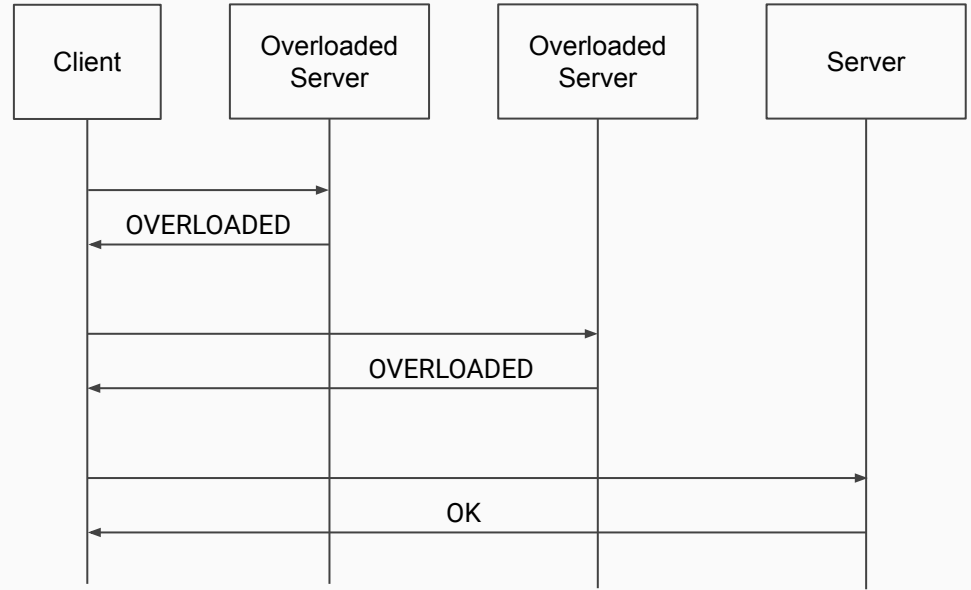
# Loadshedding

Rejecting 1 request is better than deadlining 2
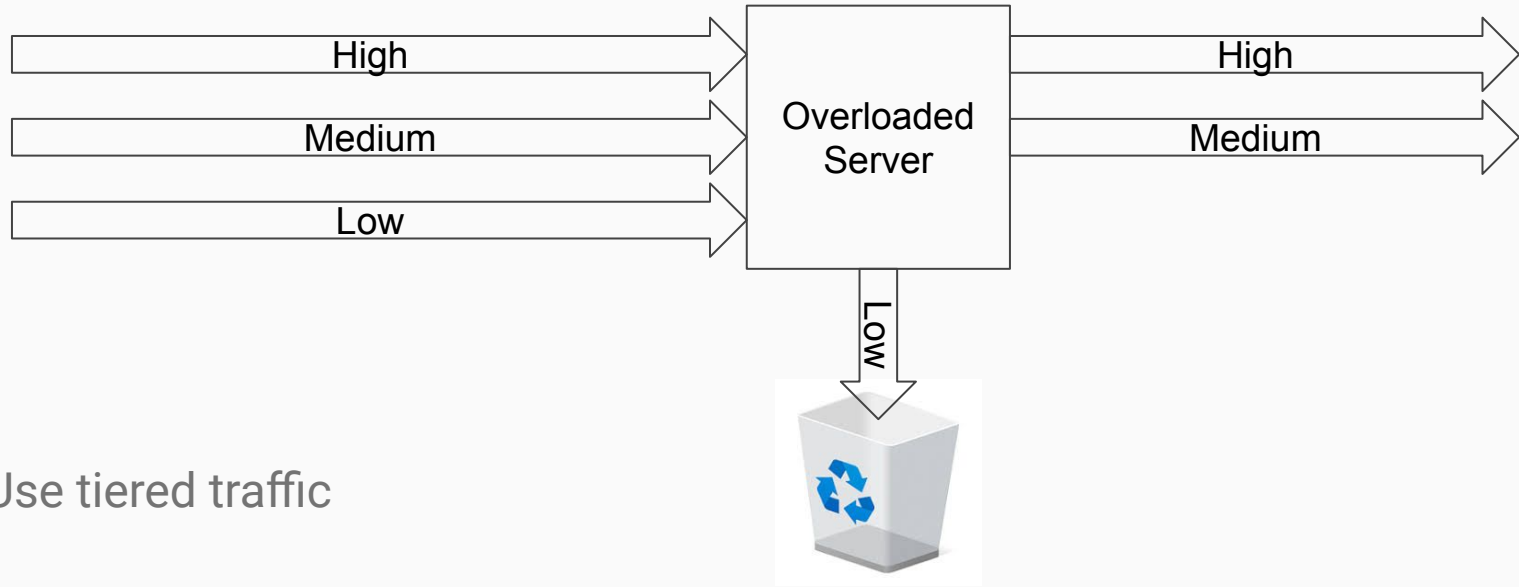
Reject fast, reject early

Can be done client side: throttle on latency/errors, exponential backoff

Loadbalancing via loadshedding

## Loadbalancing via loadshedding

# Quality Of Service



```
High ──────────────▶ ┌──────────┐ ──────────────▶ High
Medium ────────────▶ │Overloaded│ ────────────────▶ Medium
Low ───────────────▶ │  Server  │
                     └────┬─────┘
                          │ Low
                          ▼
                       [trash]
```
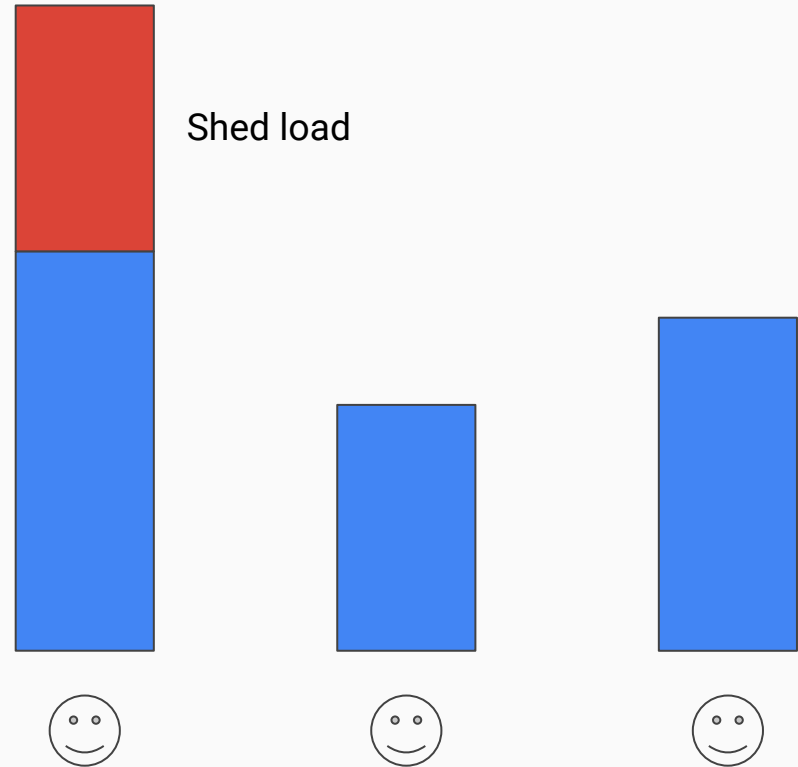
- Use tiered traffic

- Serve degraded results

# Fairness Under Overload

Maintain performance for as many
users as possible
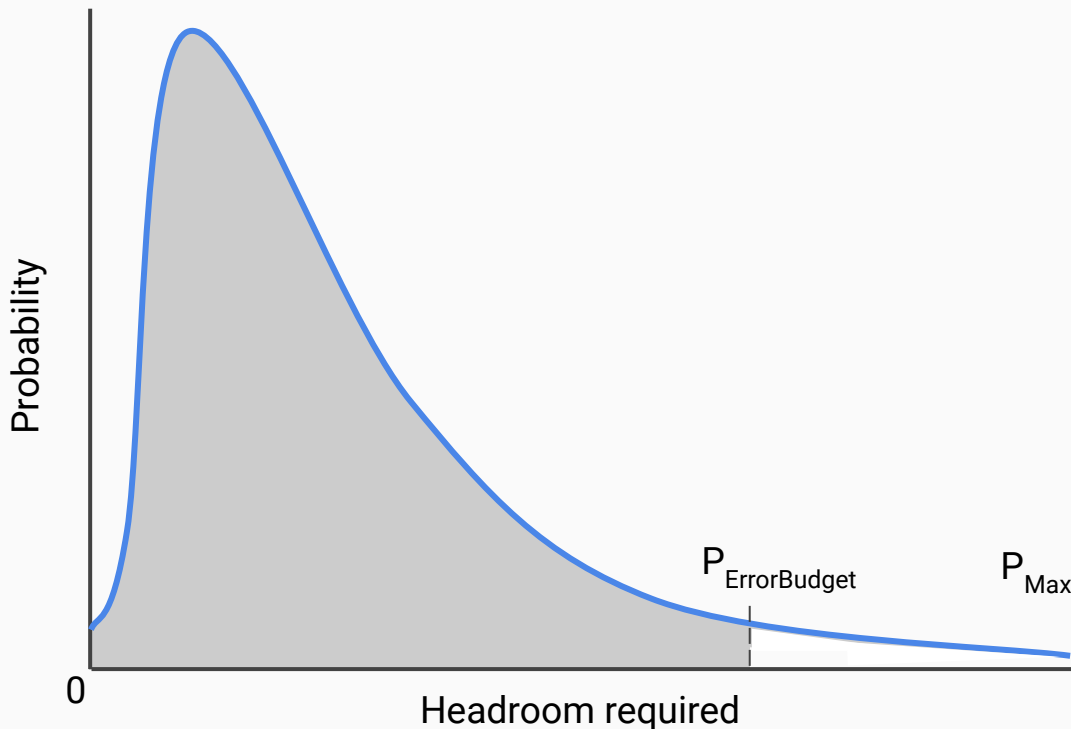
Choose the right scheduler

Load shed proportionally

Shed load

Efficiency = safe_cpu% − $P_{max}$(

   traffic_spike

   + load_imbalance$_{max-mean}$

   + failover_capacity

)

⬇

Efficiency = safe_cpu% − $\mathbf{P_{ErrorBudget}}$(

   traffic_spike

   + load_imbalance$_{max-mean}$

   + failover_capacity

)



Probability

0

Headroom required

$P_{ErrorBudget}$

$P_{Max}$

# Takeaways

$$Efficiency = safe\_cpu\% - P_{ErrorBudget}(traffic\_spike + load\_imbalance + failover\_capacity)$$

- Autoscale everything.

- Expand your footprint to minimise N+1 overheads.

- Test your failovers regularly.

- Understand your resource needs. Reshape your servers to fit. Trade off resources.

- Being more reliable under overload improves efficiency.

# Questions?

#22apac-day2-track1