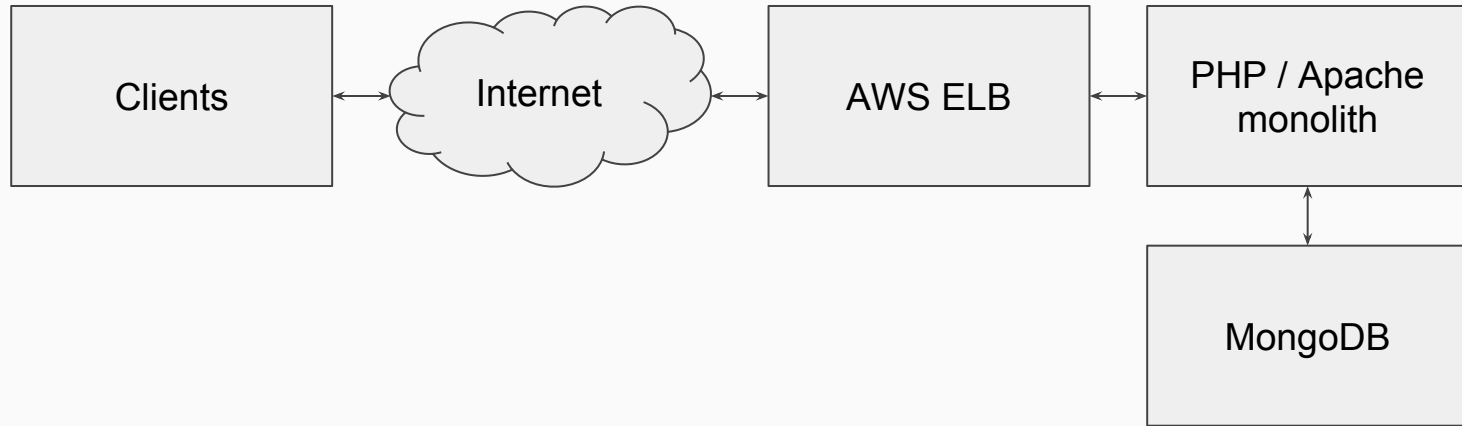Lyft's Envoy: Experiences Operating a Large Service Mesh
SREcon17
Matt Klein / @mattklein123, Software Engineer @Lyft

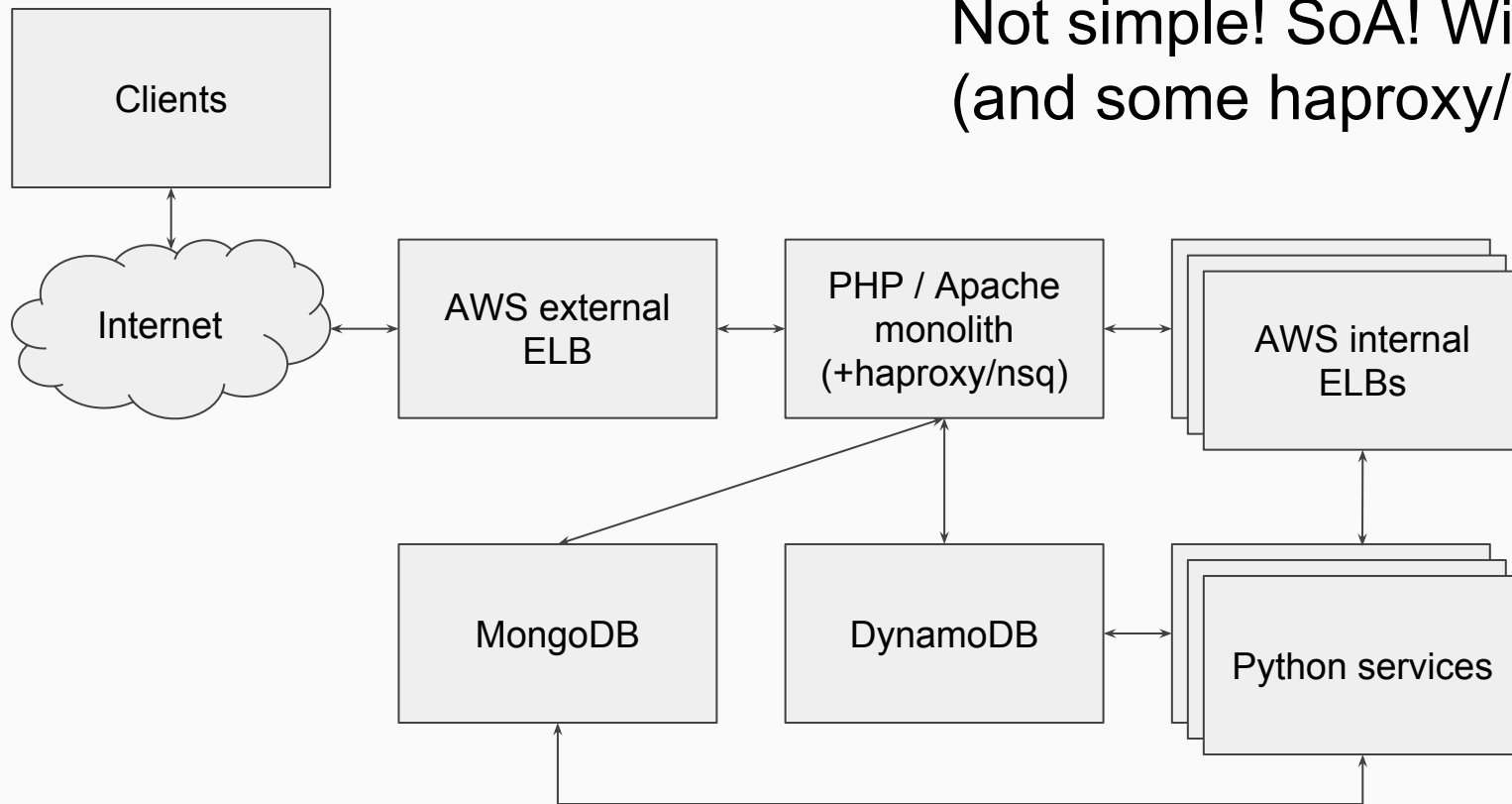Simple! No SoA! (*but still not that simple*)

# State of SoA networking in industry

- **Languages** and frameworks.
- **Protocols** (HTTP/1, HTTP/2, gRPC, databases, caching, etc.).
- **Infrastructures** (IaaS, CaaS, on premise, etc.).
- Intermediate **load balancers** (AWS ELB, F5, etc.).
- **Observability** output (stats, tracing, and logging).
- Implementations (often partial) of **retry**, **circuit breaking**, **rate limiting**, **timeouts**, and other distributed systems best practices.
- **Authentication** and **Authorization**.
- Per language **libraries** for service calls.
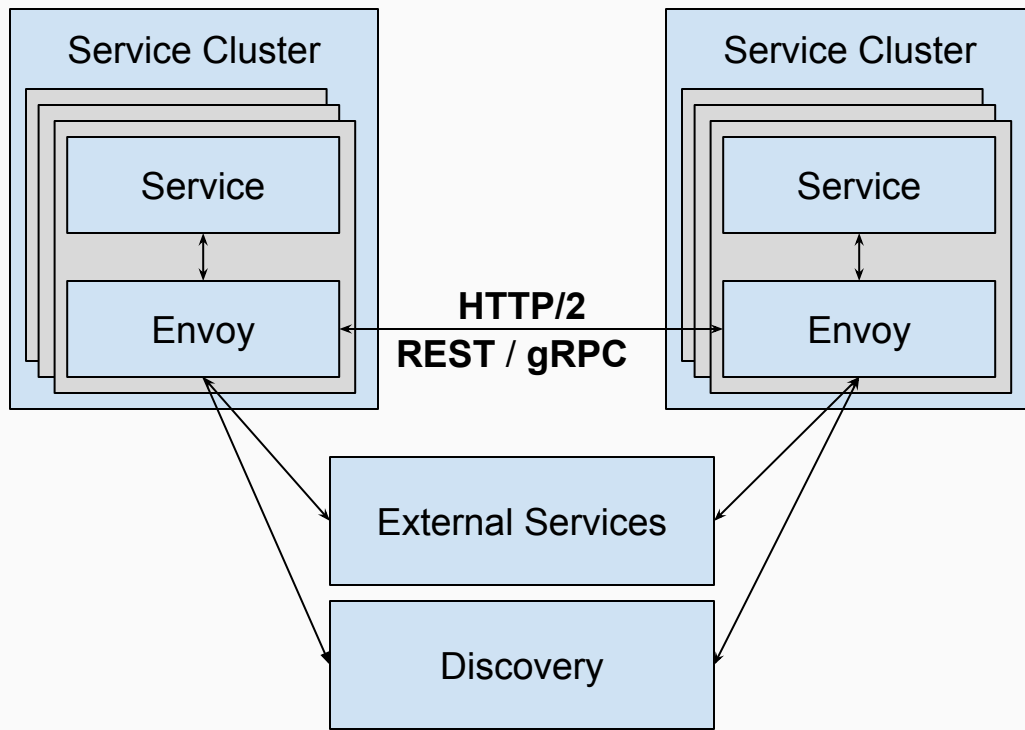
A really big and confusing mess...

*The network should be transparent to applications. When network and application problems do occur it should be easy to determine the source of the problem.*

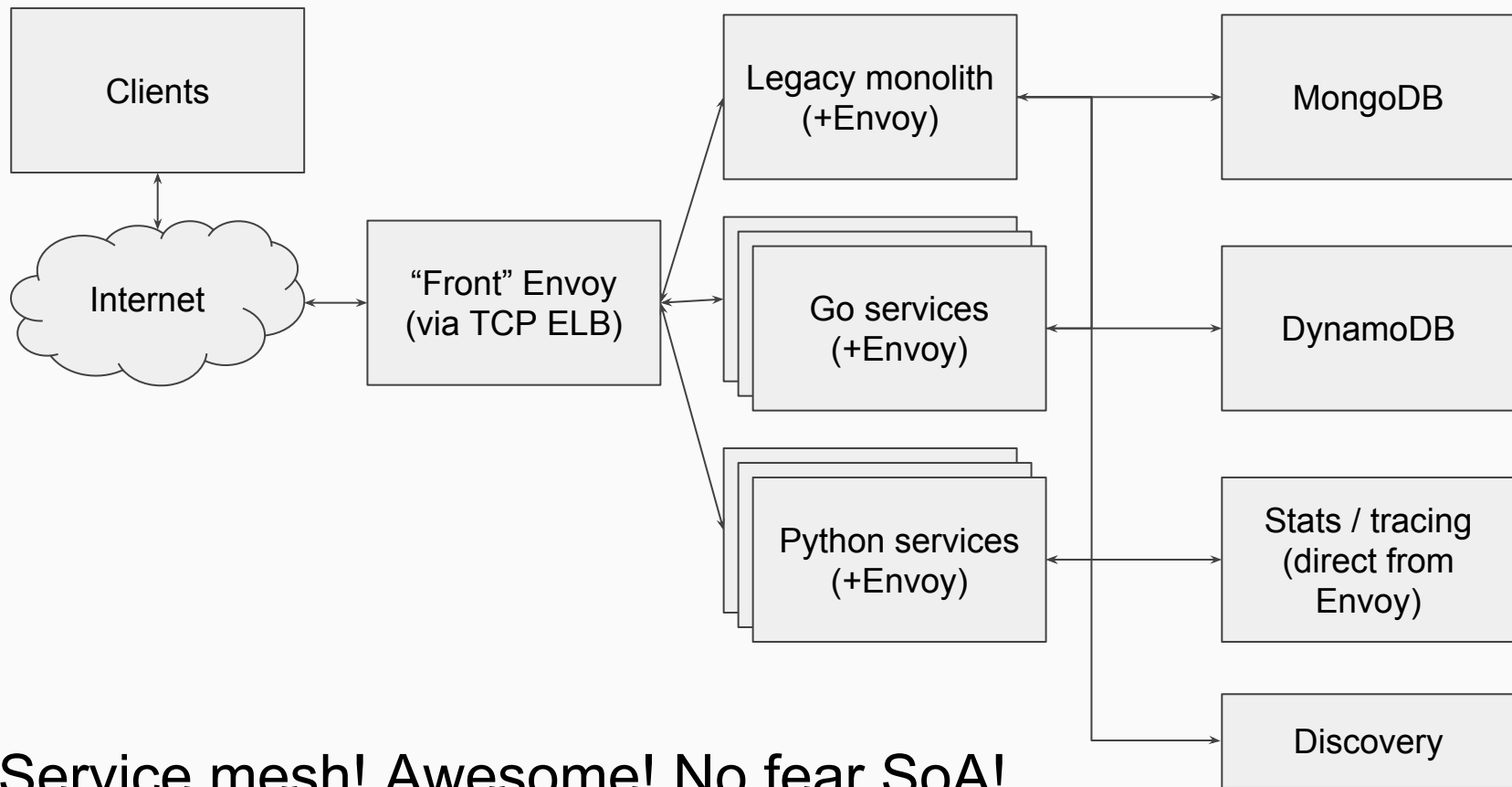This sounds great! But it turns out it's really, really hard.

# What is Envoy

- **Out of process architecture**: Let's do a lot of really hard stuff in one place and allow application developers to focus on business logic.
- **Modern C++11 code base**: Fast and productive.
- **L3/L4 filter architecture**: A byte proxy at its core. Can be used for things other than HTTP (e.g., MongoDB, redis, stunnel replacement, TCP rate limiter, etc.).
- **HTTP L7 filter architecture**: Make it easy to plug in different functionality.
- **HTTP/2** first! (Including **gRPC** and a nifty gRPC HTTP/1.1 bridge).
- **Service discovery** and **active/passive health checking.**
- **Advanced load balancing**: Retry, timeouts, circuit breaking, rate limiting, shadowing, outlier detection, etc.
- Best in class **observability**: stats, logging, and tracing.
- **Edge proxy**: routing and TLS.

# Envoy service to service topology

# Lyft today

Clients

Internet

"Front" Envoy (via TCP ELB)

Legacy monolith (+Envoy)

Go services (+Envoy)

Python services (+Envoy)

MongoDB

DynamoDB

Stats / tracing (direct from Envoy)

Discovery

## Service mesh! Awesome! No fear SoA!

# Eventually consistent service discovery

- **Fully consistent** service discovery systems are very popular (ZK, etcd, consul, etc.).
- In practice they are **hard to run at scale**.
- Service discovery is actually an **eventually consistent** problem. Let's recognize that and design for it.
- Envoy is designed from the get go to treat **service discovery as lossy**.
- Active health checking used in combination with service discovery to produce a **routable overlay**.

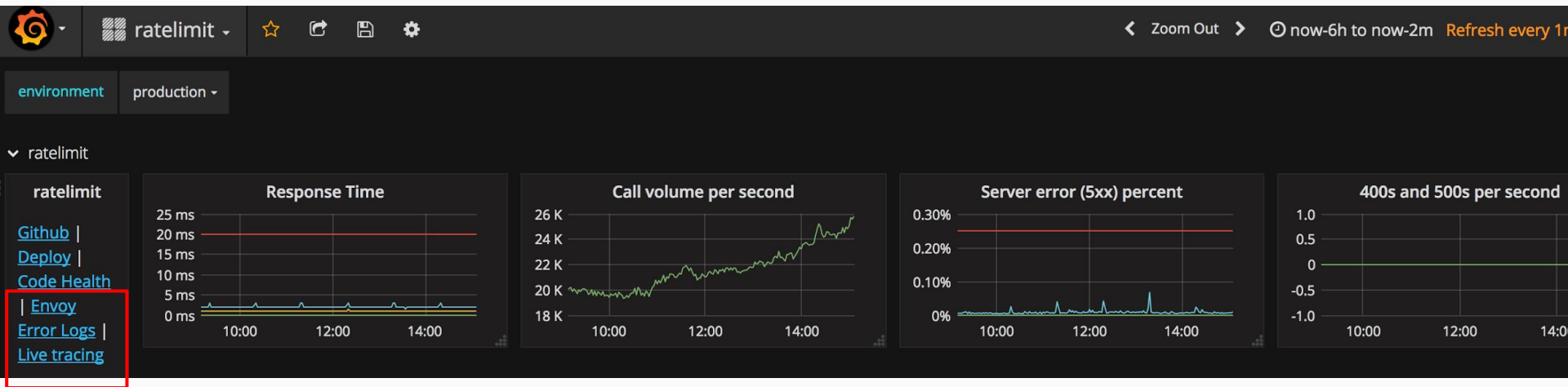| Discovery Status | HC OK | HC Failed |
|---|---|---|
| Discovered | Route | Don't Route |
| Absent | Route | Don't Route / Delete |

# Advanced load balancing

- Different **service discovery** types.
- Zone aware least request **load balancing**.
- **Dynamic stats**: Per zone, canary specific stats, etc.
- **Circuit breaking**: Max connections, requests, and retries.
- **Rate limiting**: Integration with global rate limit service.
- **Shadowing**: Fork traffic to a test cluster.
- **Retries**: HTTP router has built in retry capability with different policies.
- **Timeouts**: Both "outer" (including all retries) and "inner" (per try) timeouts.
- **Outlier detection**: Consecutive 5xx
- **Deploy control**: Blue/green, canary, etc.
- **Fault injection**

- **Observability** is by far the most important thing that Envoy provides.
- Having all SoA traffic transit through Envoy gives us a single place where we can:
  - Produce consistent **statistics** for every hop
  - Create and propagate a stable **request ID / tracing context**
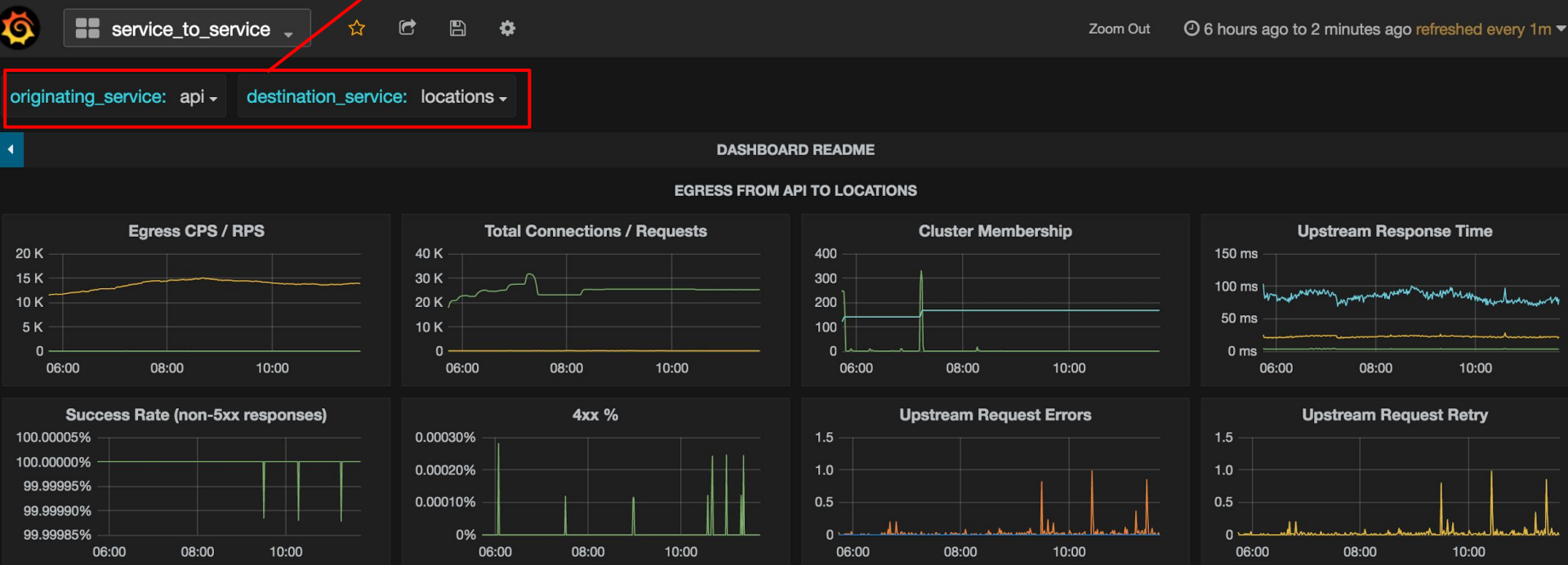  - Consistent **logging**
  - Distributed **tracing**

# Observability: Per service auto-generated panel



Links to logging and tracing

# Observability: Service to service template dashboard

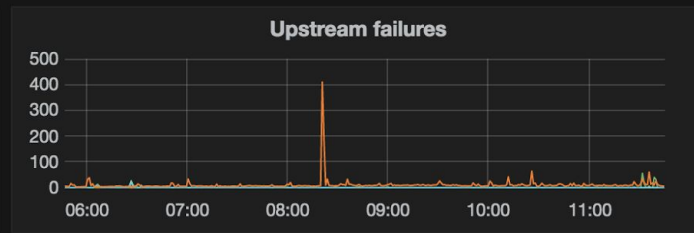Template with drop down for every service



service_to_service

Zoom Out  ⏱ 6 hours ago to 2 minutes ago refreshed every 1m ▼

originating_service: api ▾    destination_service: locations ▾

◀ DASHBOARD README

**EGRESS FROM API TO LOCATIONS**

**Egress CPS / RPS**

20 K
15 K
10 K
5 K
0
06:00    08:00    10:00

**Total Connections / Requests**

40 K
30 K
20 K
10 K
0
06:00    08:00    10:00

**Cluster Membership**

400
300
200
100
0
06:00    08:00    10:00

**Upstream Response Time**

150 ms
100 ms
50 ms
0 ms
06:00    08:00    10:00

**Success Rate (non-5xx responses)**

100.00005%
100.00000%
99.99995%
99.99990%
99.99985%
06:00    08:00    10:00

**4xx %**

0.00030%
0.00020%
0.00010%
0%
06:00    08:00    10:00

**Upstream Request Errors**

1.5
1.0
0.5
0
06:00    08:00    10:00

**Upstream Request Retry**

1.5
1.0
0.5
0
06:00    08:00    10:00

# Observability: Envoy global health dashboard

# Observability: Distributed tracing

# Observability: Logging

# Performance matters for a service proxy

- For most companies **developer time is worth more than infra costs** (cost vs. throughput).

- However, **Latency** and **predictability** is what matters. And in particular **tail latency** (P99+).

- Virtual IaaS, multiple languages and runtimes, languages that use GC: Niceties that improve productivity and reduce upfront dev costs, but make **debugging really difficult**.

- Ability to **reason about overall performance** and reliability is critical.

# Envoy config/process management @Lyft

```
Jinja JSON                 Service
templates                  manifests
    │  \                       │
    │   \                      │
    ▼    \                     ▼
"Front" Envoy   \─────▶  Service/Envoy
build/deploy                 deploy
    │                          │
    ▼                          ▼
Binaries/configs          StS Envoy
    │                      configs
     \                        /
      \                      /
       ▼                    ▼
           Salt/runit
```
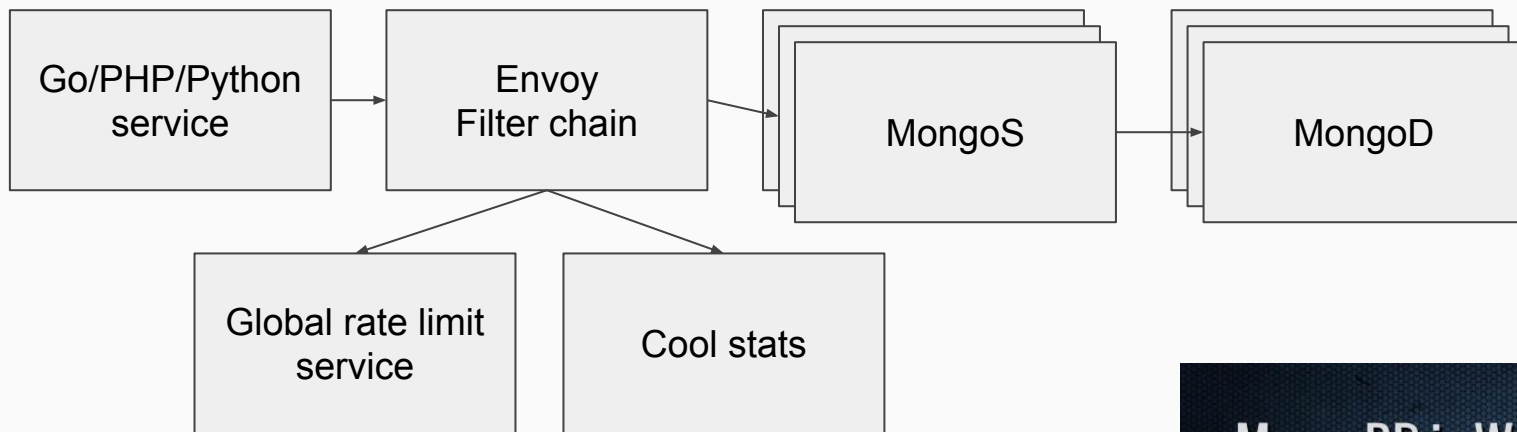
- Combination of static and dynamic configs.
- Service egress, circuit breaking, etc. configs specified in manifest.
- Service configs built on service host at service/envoy deploy time.
- **Next up for Lyft: config service via APIs!**

# Envoy thin clients @Lyft

```python
from lyft.api_client import EnvoyClient
switchboard_client = EnvoyClient(
    service='switchboard'
)
msg = {'template': 'breaksignout'}
headers = {'x-lyft-user-id': 12345647363394}
switchboard_client.post("/v2/messages", data=msg, headers=headers)
```

- Abstract away egress port
- Request ID/tracing propagation
- Guide devs into good timeout, retry, etc. policies
- Similar thin clients for Go and PHP

```
┌─────────────────┐     ┌─────────────────┐     ┌─────────────────┐     ┌─────────────────┐
│  Go/PHP/Python  │────▶│     Envoy       │────▶│                 │────▶│                 │
│    service      │     │  Filter chain   │     │     MongoS      │     │     MongoD      │
└─────────────────┘     └─────────────────┘     └─────────────────┘     └─────────────────┘
                               │     │
                        ┌──────┘     └──────┐
                        ▼                   ▼
               ┌─────────────────┐ ┌─────────────────┐
               │ Global rate limit│ │                 │
               │     service      │ │    Cool stats   │
               └─────────────────┘ └─────────────────┘
```

**No more death spirals! Web scale!**

Filters:
- L4 global rate limit (limit CPS into MongoS)
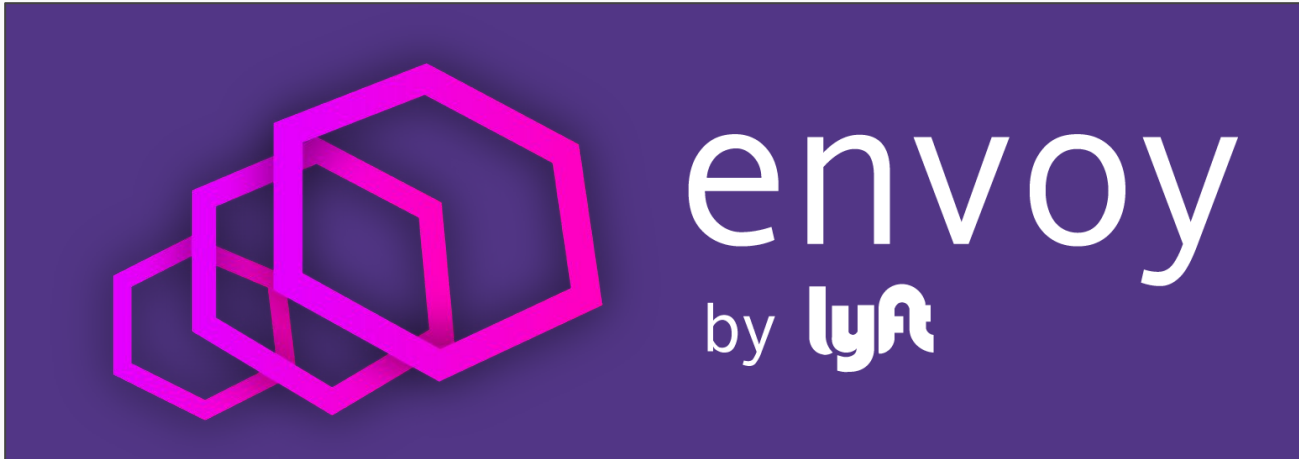- Mongo sniffer (cool stats)
- TCP proxy (MongoS load balancing)



MongoDB is Web Scale

- > 100 services.
- > 10,000 hosts.
- > 2,000,000 RPS.
- All service to service traffic (REST and gRPC).
- Use gRPC bridge to unlock Python gevent clients.
- MongoDB proxy.
- DynamoDB proxy.
- External service proxy (AWS and other partners).
- Kibana/Elastic Search for logging.
- LightStep for tracing.
- Wavefront for stats (via statsd).

- Thanks for coming!
- We are super excited about building a **community** around Envoy. Talk to us if you need help getting started.
- https://lyft.github.io/envoy/
- **Lyft is hiring**: Contact us if you want to work on hard scaling problems in a fast moving company: https://www.lyft.com/jobs