# Resilient Routing and Discovery

**Simon Eskildsen, Shopify**
**@Sirupsen**

dockercon 15

SF JUNE 22-23

# Shopify

## Docker in Production serving the below for 1+ year

**165,000+**

ACTIVE SHOPIFY MERCHANTS

**200+**

DEVELOPERS

**2**

DATACENTERS

**3000+**

CONTAINERS RUNNING AT ANY TIME

**12+**

DEPLOYS PER DAY

**10,000+**

MAX CHECKOUTS PER MINUTE

**$8 BILLION+**

CUMULATIVE GMV

**500+**

SERVERS

**Ruby on Rails**

10+ years old

**300M unique visits/month**

LEAGUE OF APPLE, EBAY AND AMAZON

**Building reliable bridges in large distributed systems**

Complexity (y-axis) vs Reliability (x-axis)

Cross DC Networking

Cross Regional Networking

Same Rack Networking

Inter process

In process

5

**Resiliency**

**Discovery**

**Routing**

**Reliability** is your success metric for discovery and routing.

**Shopify started this journey in the fall of 2014**

# Resiliency

Building a reliable system from unreliable components

# (Micro)service equation

$$Uptime = A^N$$

Number of services

Availability per service

Total availability

# Resiliency Matrix

|  | Checkout | Admin | Storefront |
|---|---|---|---|
| MySQL Shard | Unavailable | Unavailable | Degraded |
| MySQL Master | Available | Unavailable | Available |
| Kafka | Available | Degraded | Available |
| External HTTP API | Degraded | Available | Unavailable |
| redis-sessions | Unavailable | Unavailable | Degraded |

# Objectives for large distributed systems

Building reliable systems from unreliable components

Explore resiliency, service discovery, routing, orchestration and the relationship between them

Recognizing and avoiding premature optimizations and overcompensation

**Application should be designed to handle fallbacks**

# KITH

SHOP ⌄  BRANDS ⌄  BLOG  LOOKBOOKS  ABOUT  LOCATIONS  SUBSCRIBE  SUPPORT ⌄

CART (1) 🛒
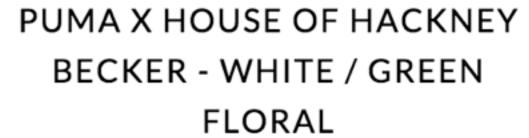
| ALL | SNEAKERS | BOOTS | SHOES | SANDALS |

NEWEST TO OLDEST ▼



### PUMA X HOUSE OF HACKNEY BECKER - WHITE / GREEN FLORAL

PUMA

$160.00



### PUMA X HOUSE OF HACKNEY BASKET - TOTAL ECLIPSE / GREEN PALM

PUMA

$125.00



### NEW BALANCE M998 "MONTAUK" - GREY / BLUE / GREEN

NEW BALANCE

$170.00



### NEW BALANCE M997 "DISTINCT" - OLIVE / BROWN

NEW BALANCE

$260.00

# Avoid HTTP 500 for single service failing

### .. or suffer the faith of the (micro)service equation



WE'LL BE *Back* SOON!

⚠️ Due to an unexpected technical problem, kithnyc.com is temporarily unavailable. Please check back in a few minutes – we'll be up and running in no time!

# Sessions data store unavailable



**Customer signed out**

KITH

SHOP ⌄      BRANDS ⌄      BLOG      LOOKBOOKS      ABOUT      LOCATIONS      SUBSCRIBE      SUPPORT ⌄      CART 🛒

Search 🔍    Account 👤

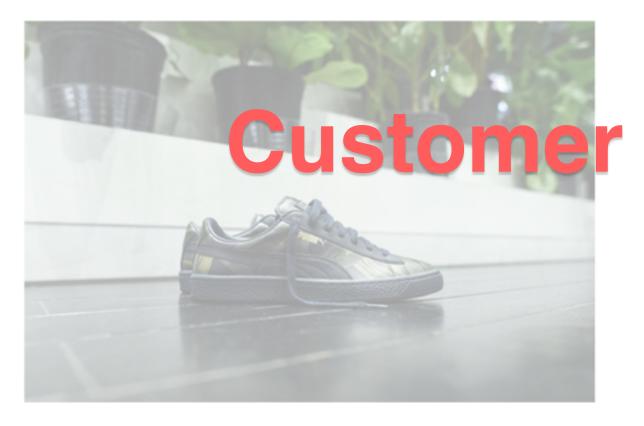ALL    SNEAKERS    BOOTS    SHOES    SANDALS

NEWEST TO OLDEST ▼

PUMA X HOUSE OF HACKNEY BECKER - WHITE / GREEN FLORAL

PUMA

$160.00

PUMA X HOUSE OF HACKNEY BASKET - TOTAL ECLIPSE / GREEN PALM

PUMA

$125.00

NEW BALANCE M998 "MONTAUK" - GREY / BLUE / GREEN

NEW BALANCE

$170.00

NEW BALANCE M997 "DISTINCT" - OLIVE / BROWN

NEW BALANCE

$260.00

# Simulate TCP conditions with Toxiproxy

```
curl -i -d '{"enabled":true, "latency":1000}' \
  localhost:8474/proxies/redis/downstream/toxics/latency

curl -i -X DELETE localhost:8474/proxies/redis

Toxiproxy[:mysql_master].downstream(:latency, latency: 1000).apply do
  Shop.first # this takes at least 1s
end

Toxiproxy[/redis/].down do
  session[:user_id] # this will throw an exception
end
```
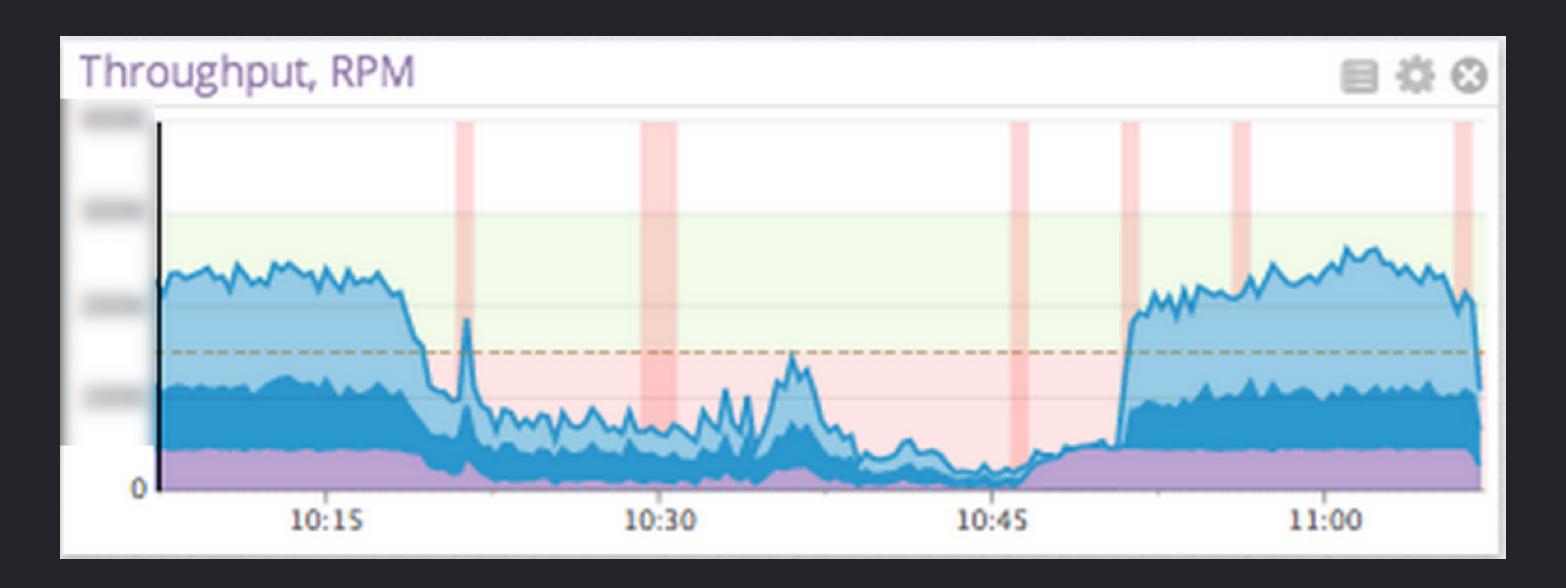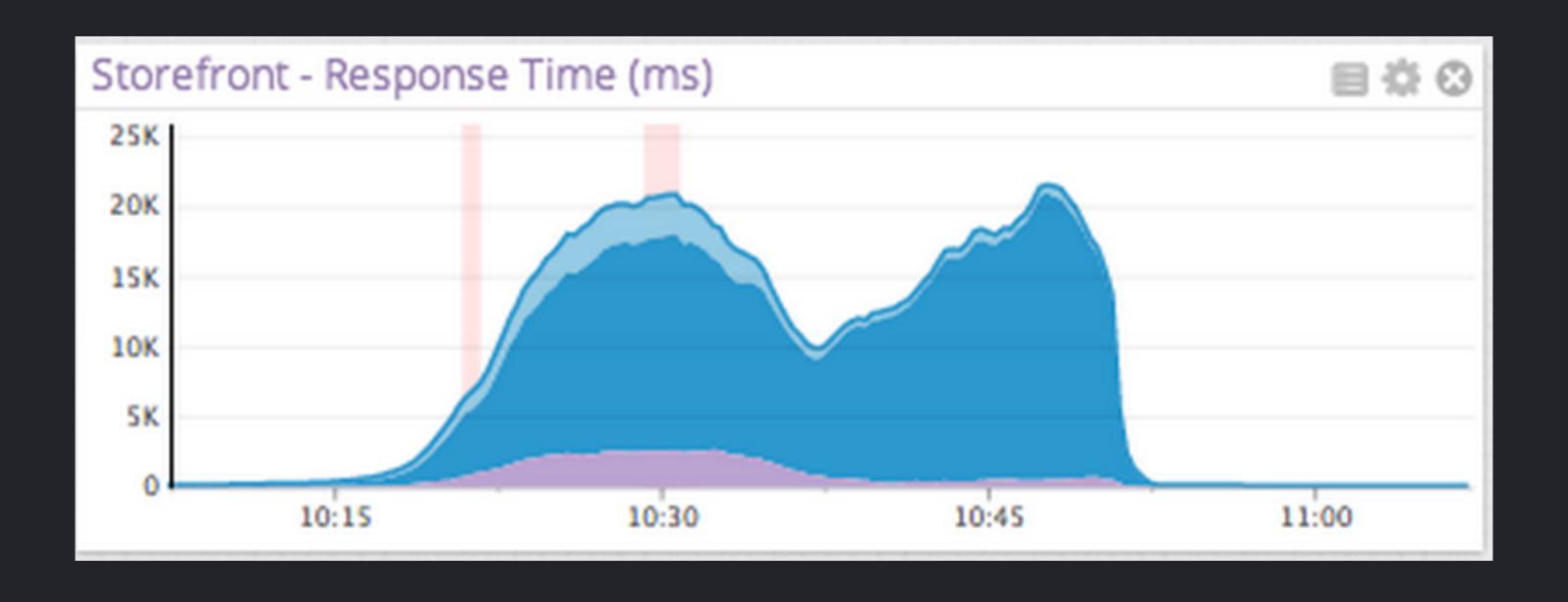
**https://github.com/shopify/toxiproxy**
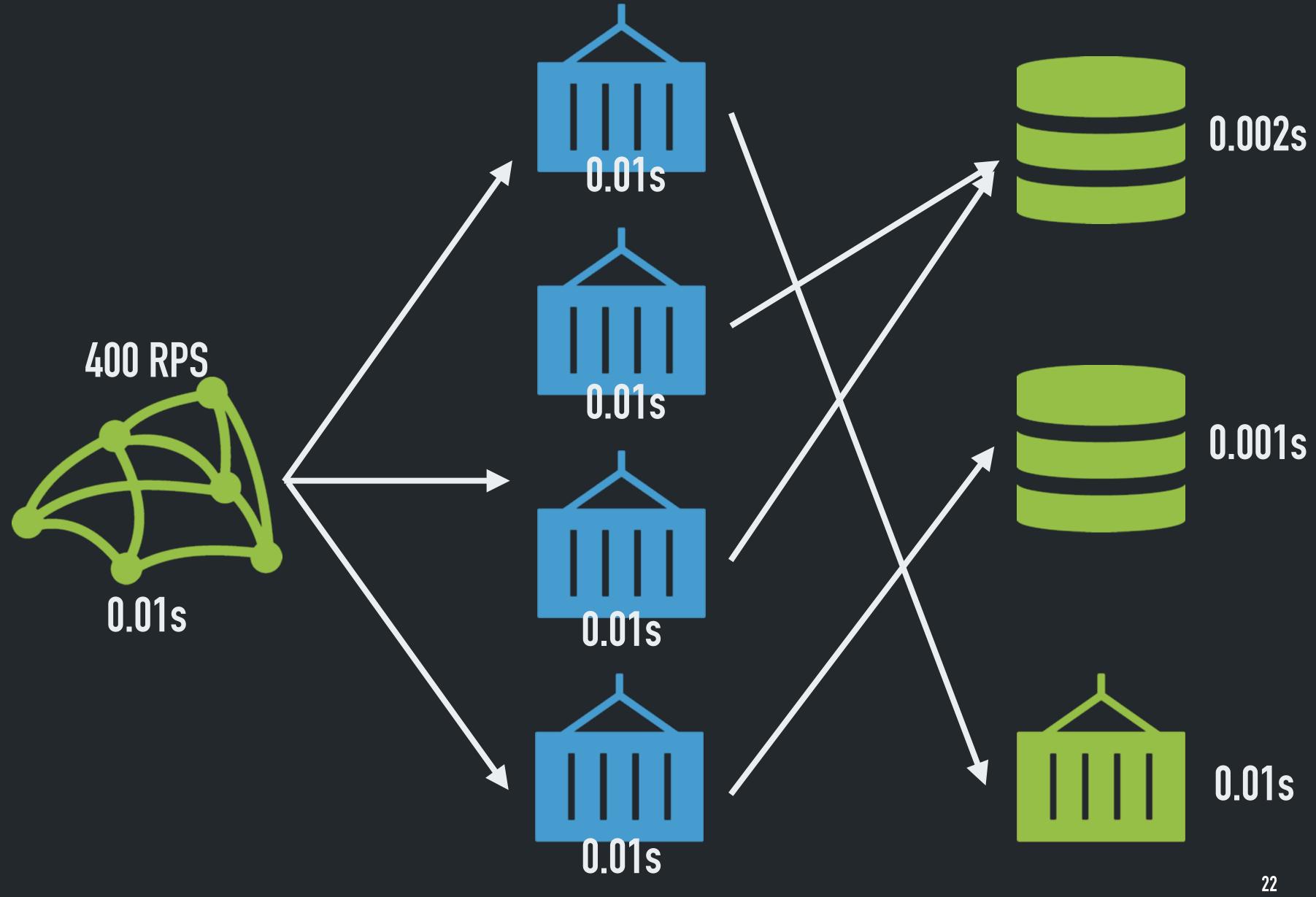
With fallbacks the system is still vulnerable to slowness. ECONNREFUSED is a **luxury, slowness is the killer.**
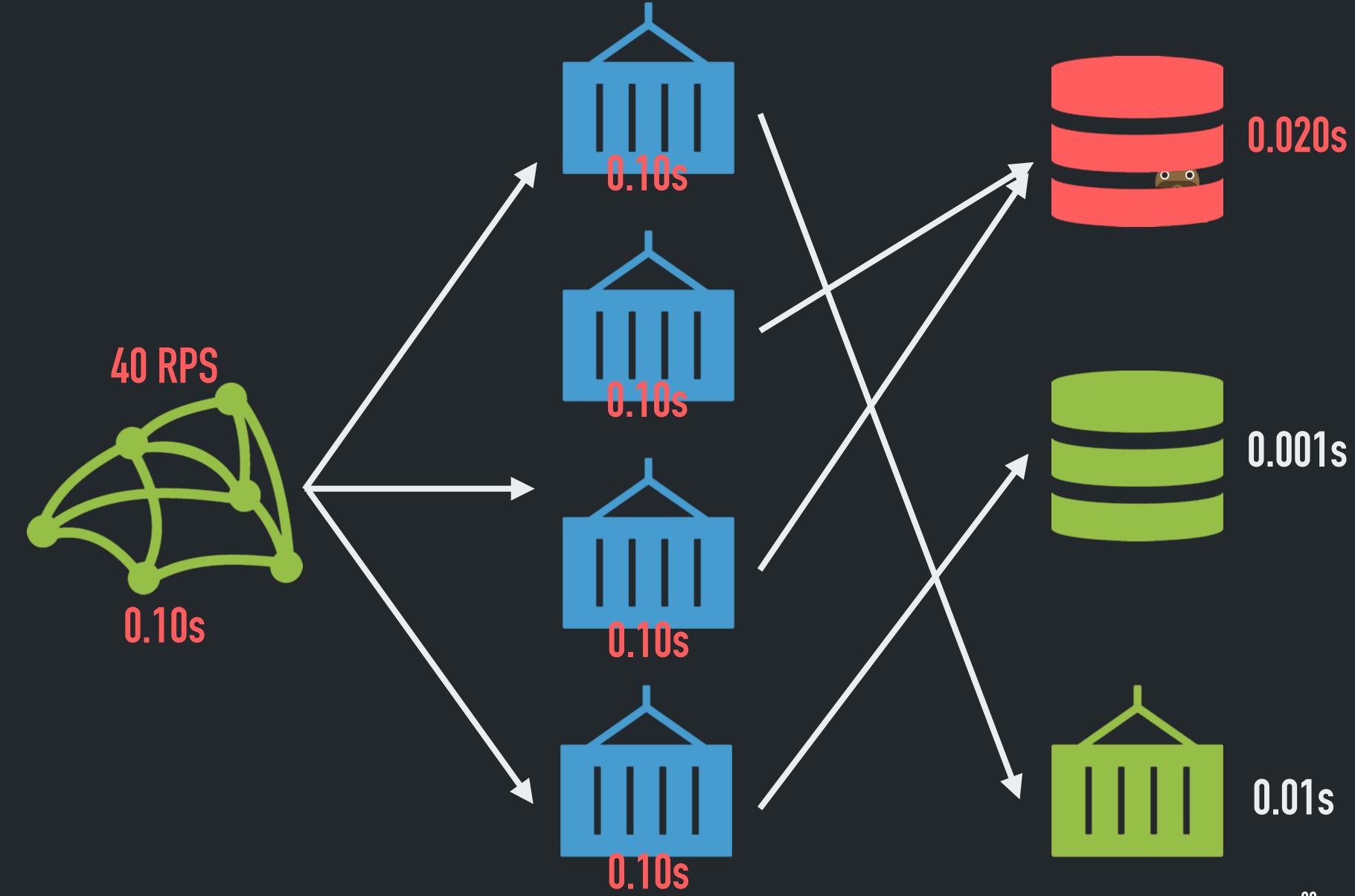
# Little's law

# Infrastructure operating normally



400 RPS
0.01s

0.01s
0.01s
0.01s
0.01s

0.002s
0.001s
0.01s

# Database latency increases by 10x, throughput drops 10x



40 RPS

0.10s

0.10s

0.10s

0.10s

0.10s

0.020s

0.001s

0.01s

**Beating Little's law is your first priority as you add services**

# Resiliency Toolkits

## Bulk Heads, Circuit Breakers, ..



Release It book


twitter/finagle


shopify/semian


netflix/hystrix

# Resiliency Maturity Pyramid

Region Gorilla

Latency Monkey

Kill Nodes (Chaos Monkey)

Production Practise Days (Games)

Resiliency Patterns

Application-Specific Fallbacks

Toxiproxy tests and matrix

Testing with mocks

No resiliency effort

# Discovery

# Infrastructure source of truth

## Services

Instances of services

## Metadata

Deployed revision, leader, ..

## Orchestration

Aid to make things happen across components

# Location



**Global**

Geo-replicated discovery

**Regional**

Single datacenter

# Discovery Backbone Properties

No single point of failure

Stale reads better than no reads: A > C

Reads order of magnitude larger than writes

Fast convergence

# New and Old School

| | |
|---|---|
| Consul | DNS |
| Zookeeper | Chef, Puppet, .. |
| Eureka | Network |
| Etcd | Hardcoded values |

**Pure DNS** for as long as you can.
Still works for us. Don't **overcompensate.**

# Pure DNS

Resilient

Simple

API

Supported

Failovers?

Slow convergence

Not a data store

Not for orchestration

**Global discovery** and **orchestration** most pressing issue for Shopify

# **Orchestration** of datacenter failovers
## **Too many Sources of Truth**

| Component | Source of Truth |
|---|---|
| Network | NetEng? |
| MySQL | DBAs? |
| Application | Cookbooks |
| Redis | Cookbooks |
| Load Balancers | Hardcode value in config file |

# Routing shops to the right datacenter



DNS: shop.walrustoys.com

CNAME
walrustoys.myshopify.com

IPs for DC 2

Map shop to DC

# DNS problematic when..

Multiple owners of data

Fast converge

Lots of change in instances

# Zookeeper

Scalable stale reads

Consistent

Orchestration

Trusted

Not complete discovery

Complex clients

Operational burden

Shoehorn

# Complex client problem

Connecting directly risky

Proxy pattern

Dumping to files

Stale reads

# Routing

# Routing responsibilities

Load balance

Protect applications against unhealthy resources: circuit breaker, bulk heads, rate limiting, …

Receive upstreams from discovery layer

| | Trusted | Scriptable | Resiliency | Dynamic upstreams | Discovery built in | TCP | Library/Proxy |
|---|---|---|---|---|---|---|---|
| **yours** | Don't do this | Of course | It's perfect | I got it | Easy | Obviously, it's Go | |
| **OS nginx** | YES | 3rd party (ngx-lua). Not complete (no TCP support). | Possible for HTTP via ngx-lua. No TCP yet | Sidekick for new upstreams. Manipulate existing via ngx-lua | No, try via sidekick/ ngx-lua | Landed in 1.9.0, stabilized in nginx+ | Proxy |
| **haproxy** | YES | Lua support in master | Not scriptable, only rate limiting built-in | Sidekick and reloads (with iptables wizardry), manipulate existing admin socket | No, try via sidekick | Built as L4 | Proxy |
| **vulcand** | Maybe? | middlewares, requires forking | SOME, only circuit breaker | Beautiful HTTP API | etcd support | No, only supports HTTP currently (not in ROADMAP.md) | Proxy |
| **finagle** | YES | YES, completely centered around plugins | YES, sophisticated FailFast module | YES | Zookeeper support | Application-level | Library, requires JVM |
| **smartstack** | Somewhat | However much HAProxy is, adapters | NO, same as HAProxy | YES | Zookeeper support | Yes, uses HAProxy | Proxy + discovery |

With a polyglot stack, we just use simple proxies and DNS

# Current Stack

## Discovery

DNS     Chef     Zookeeper

## Server

ZK Proxy

## Discoverable

## Through proxy
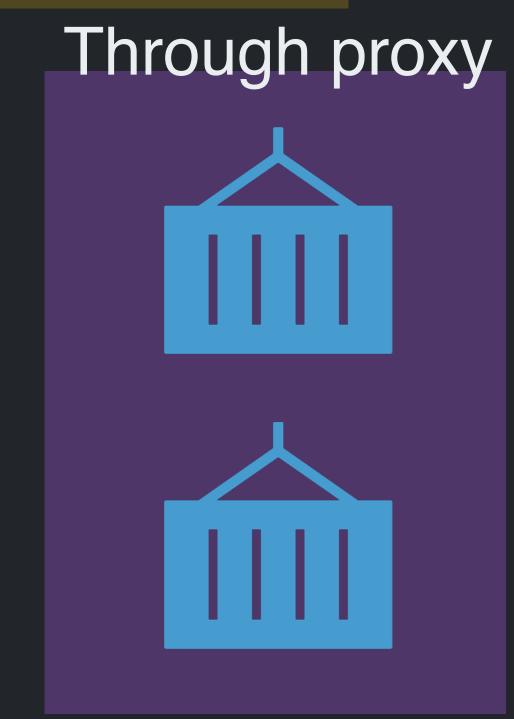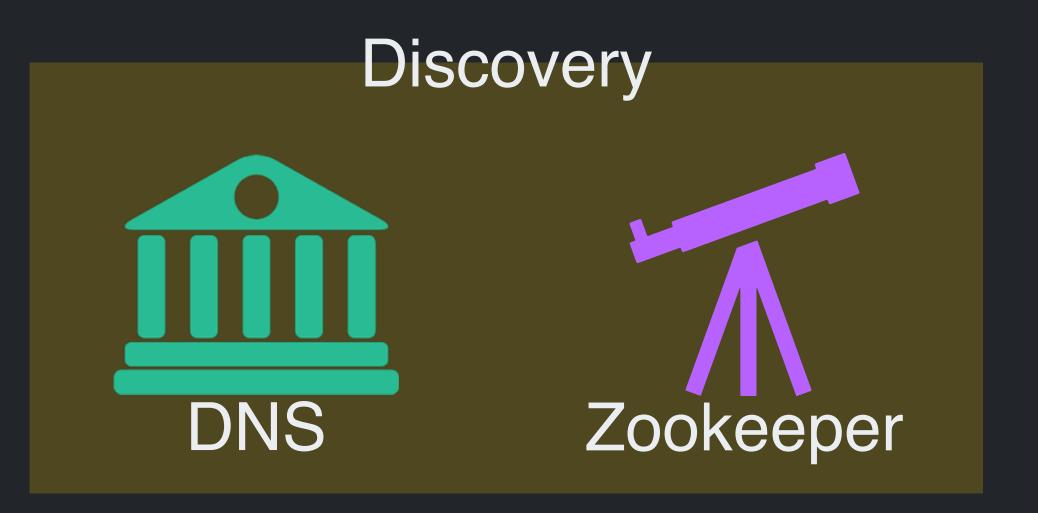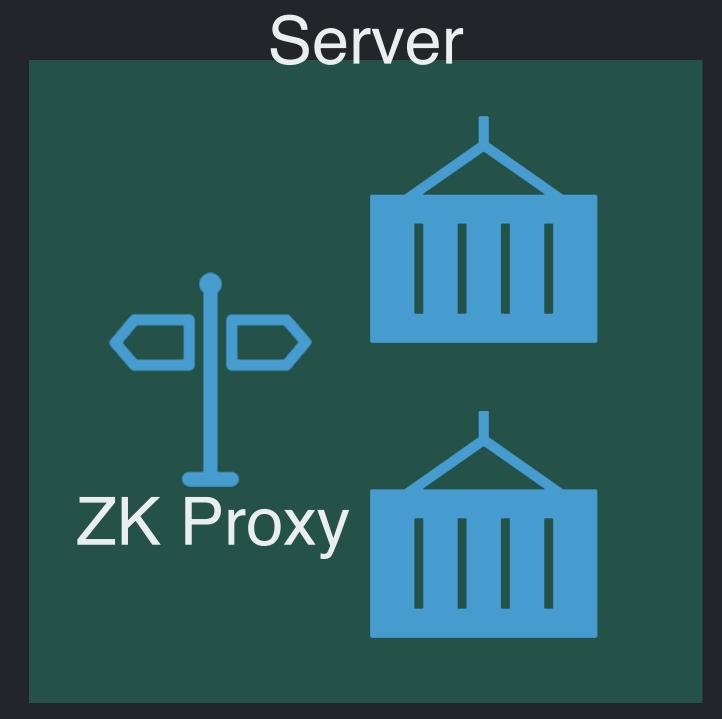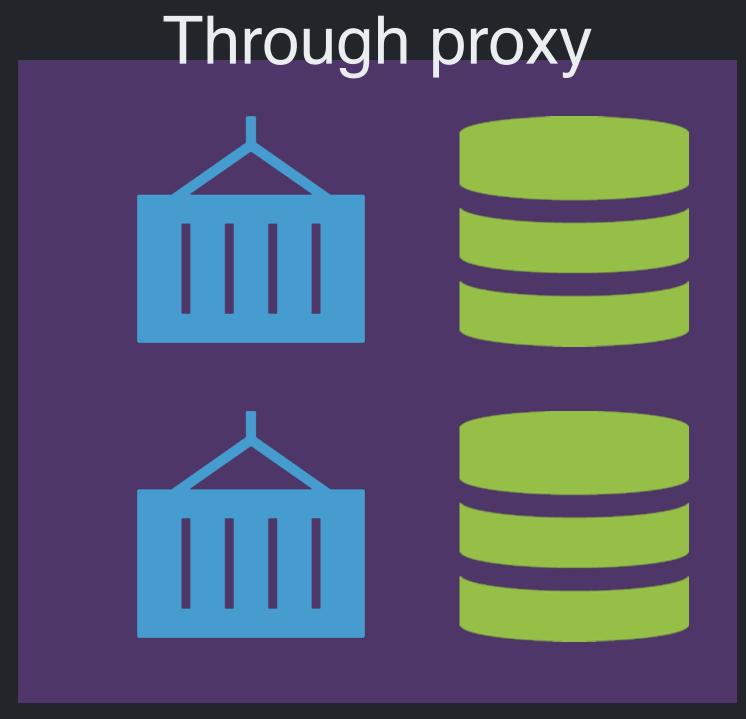
# Future Stack

# Docker's future role in discovery, routing and resiliency

# Final remarks

Figure out service discovery value for your company, don't overcompensate—your metric is reliability

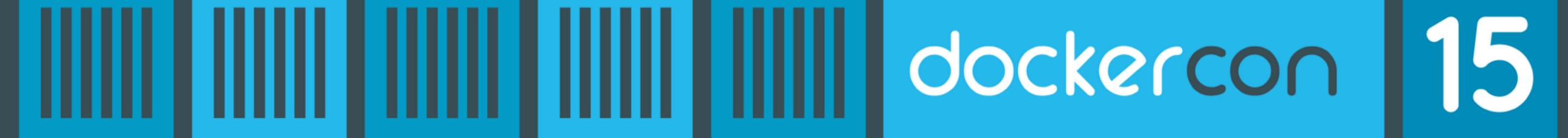Infrastructure teams own integration points, don't leave it up to everyone to jump in

Build resiliency into the system, don't make it opt in, be able to reason about entire system's state and test

# Thank you

Simon Eskildsen, Shopify

@Sirupsen

dockercon 15