

# Introducing Envoy-based service mesh at Booking.com

**Ivan Kruglov**

KubeCon Europe 2018

02.05.2018

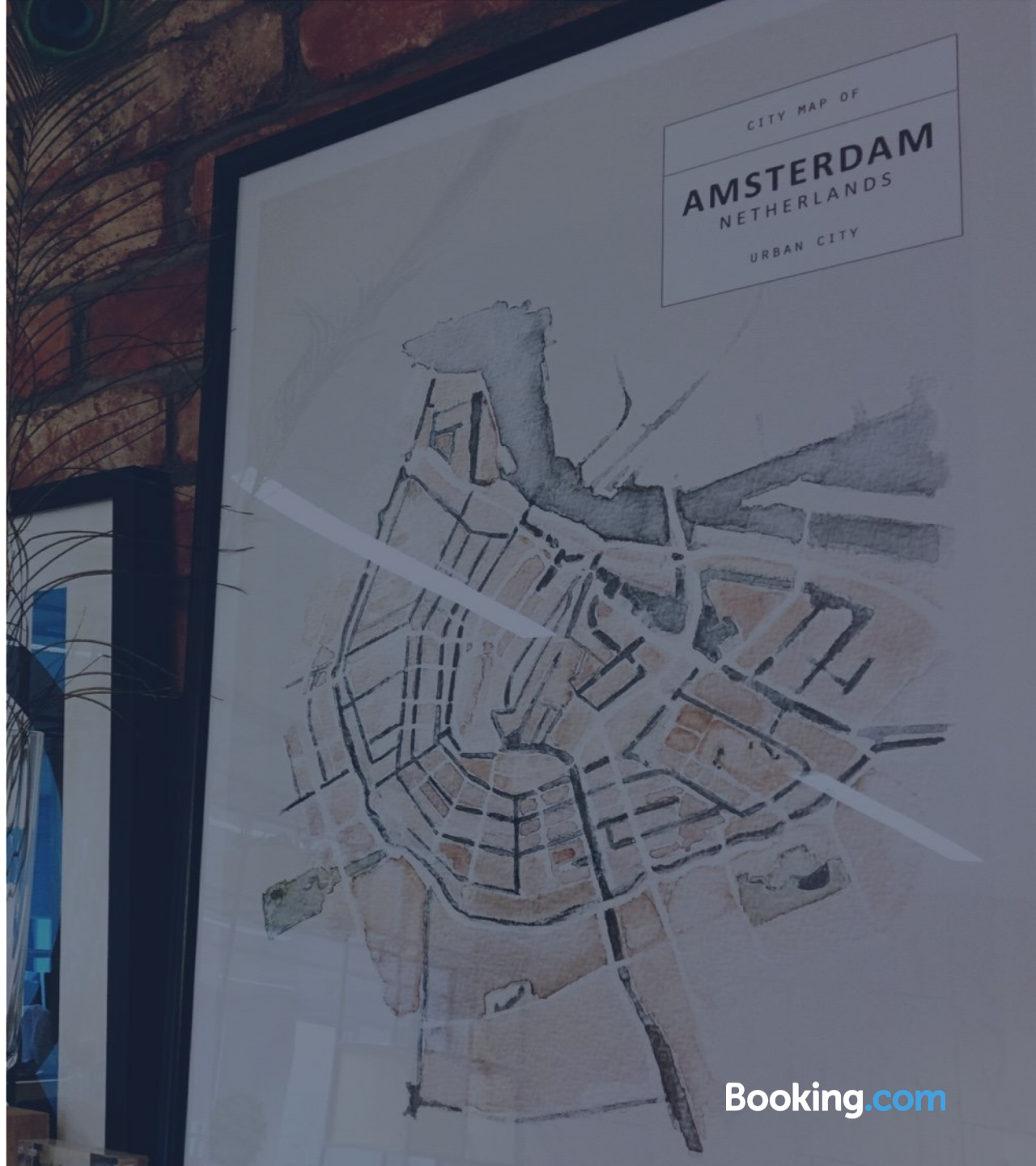
# Booking.com

based in Amsterdam

28M listings

1,5M room nights per day

1700 IT staff

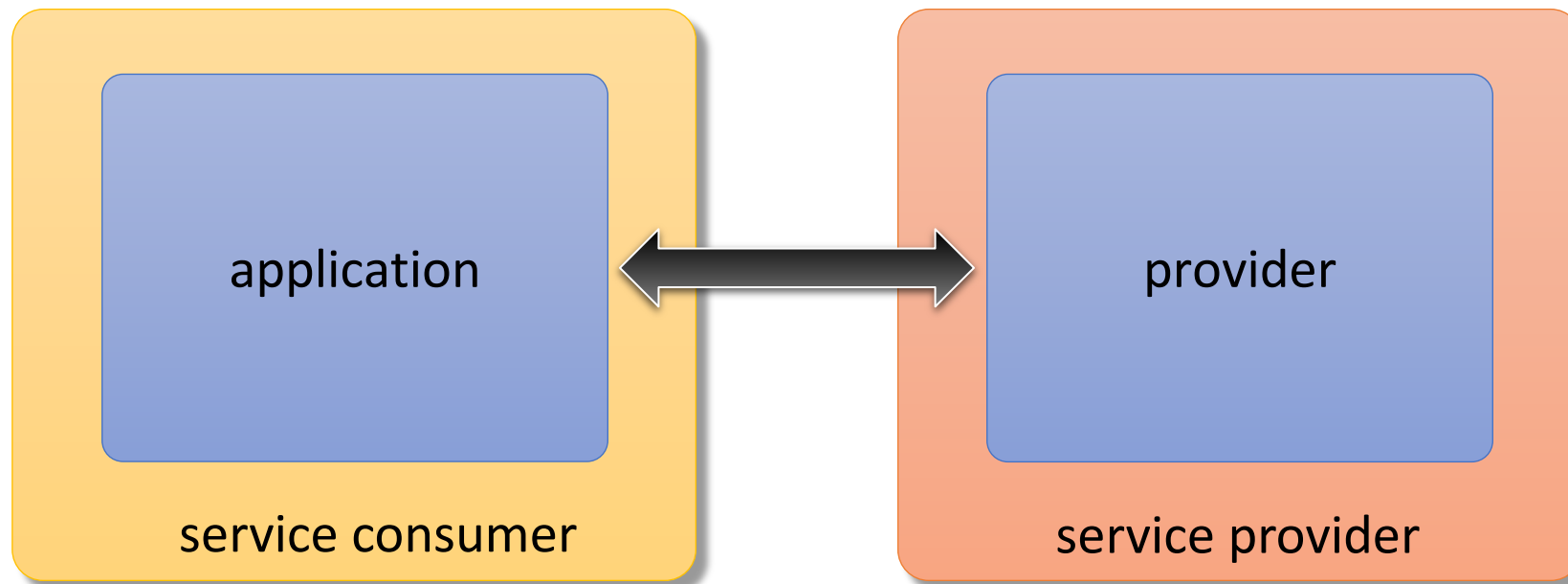


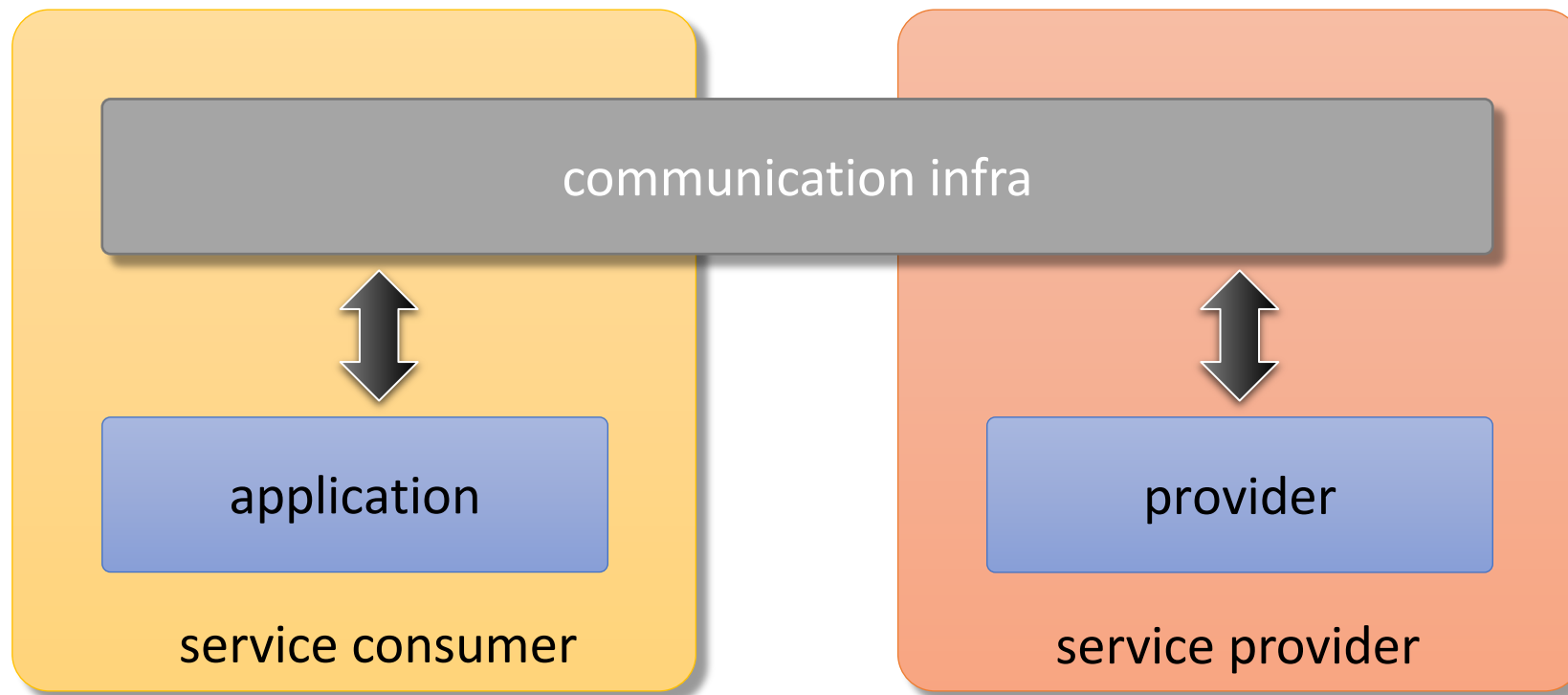
# Agenda

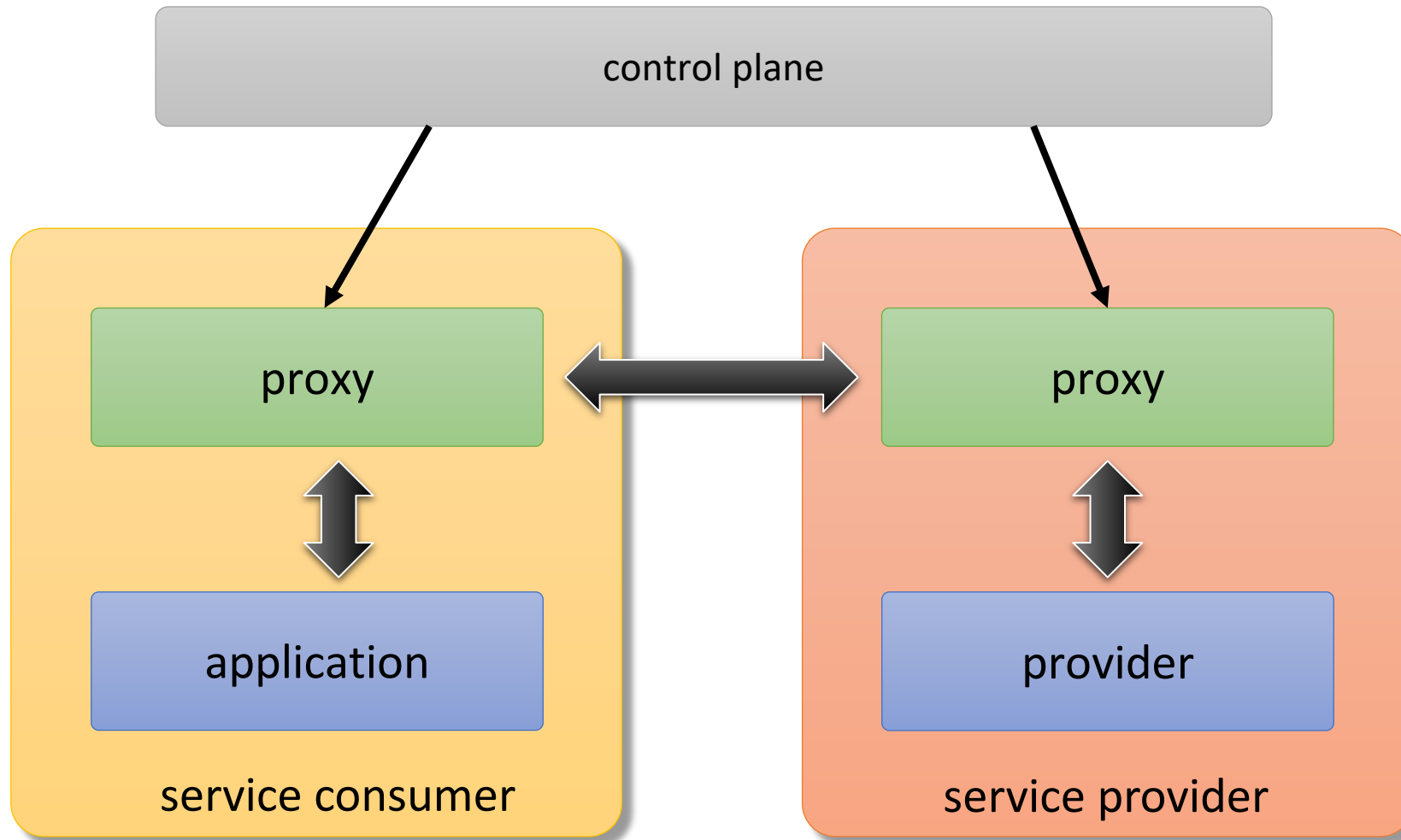
- what is service mesh?
- why did we start the project?
- our setup
- learnings
- conclusion

# What is service mesh\*?

\* the way I understand it





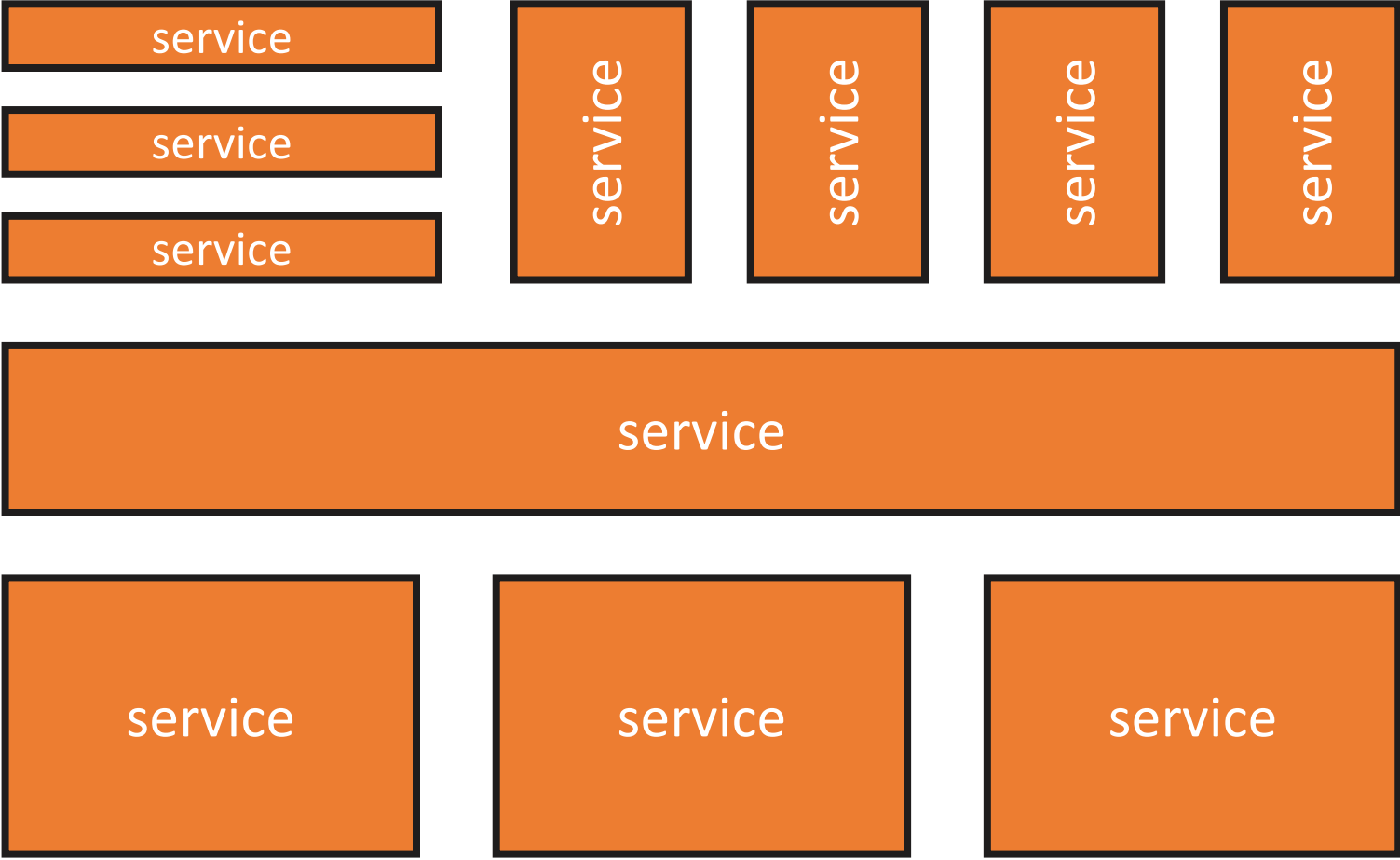


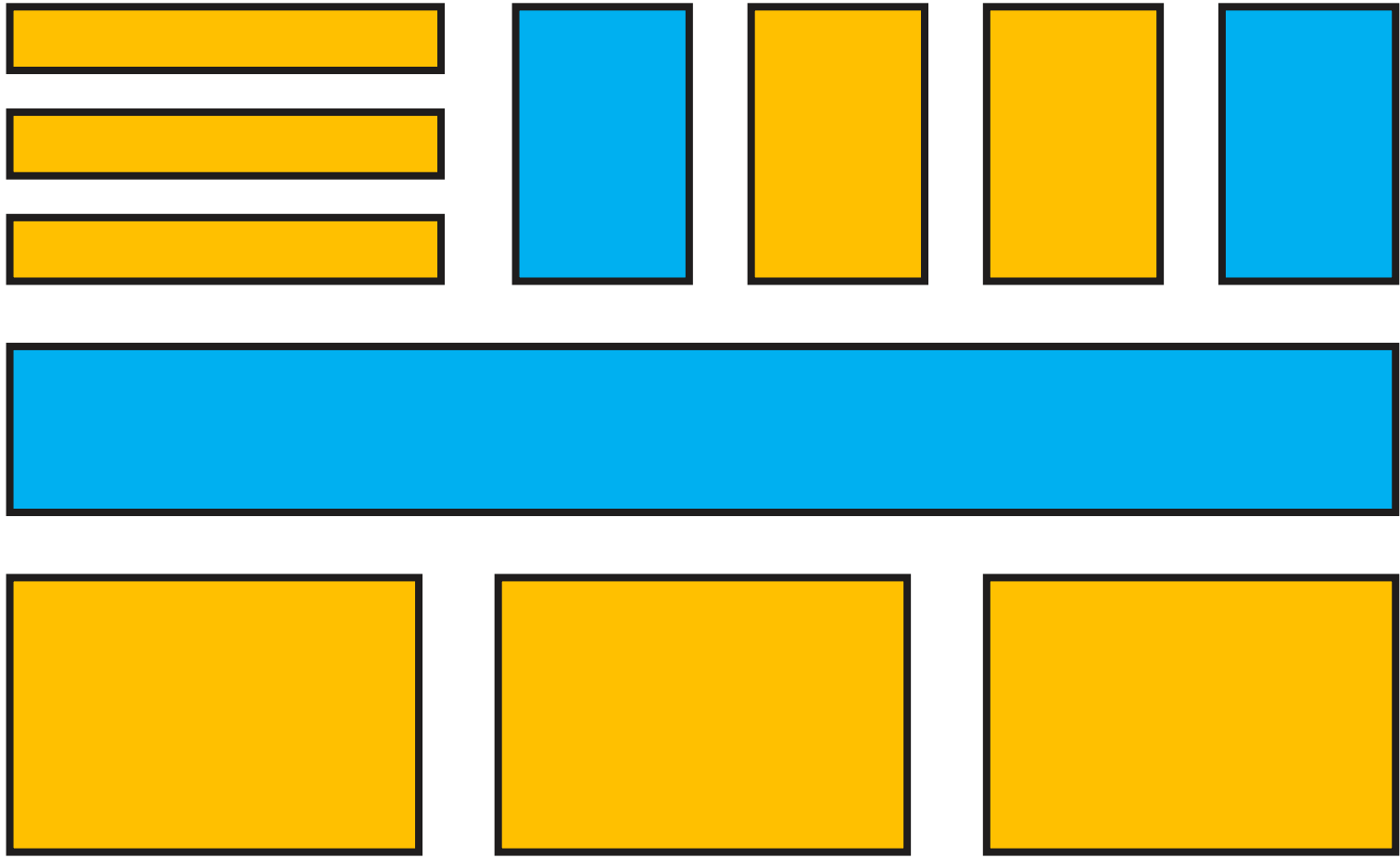
**Why did we start  
service mesh  
project?**

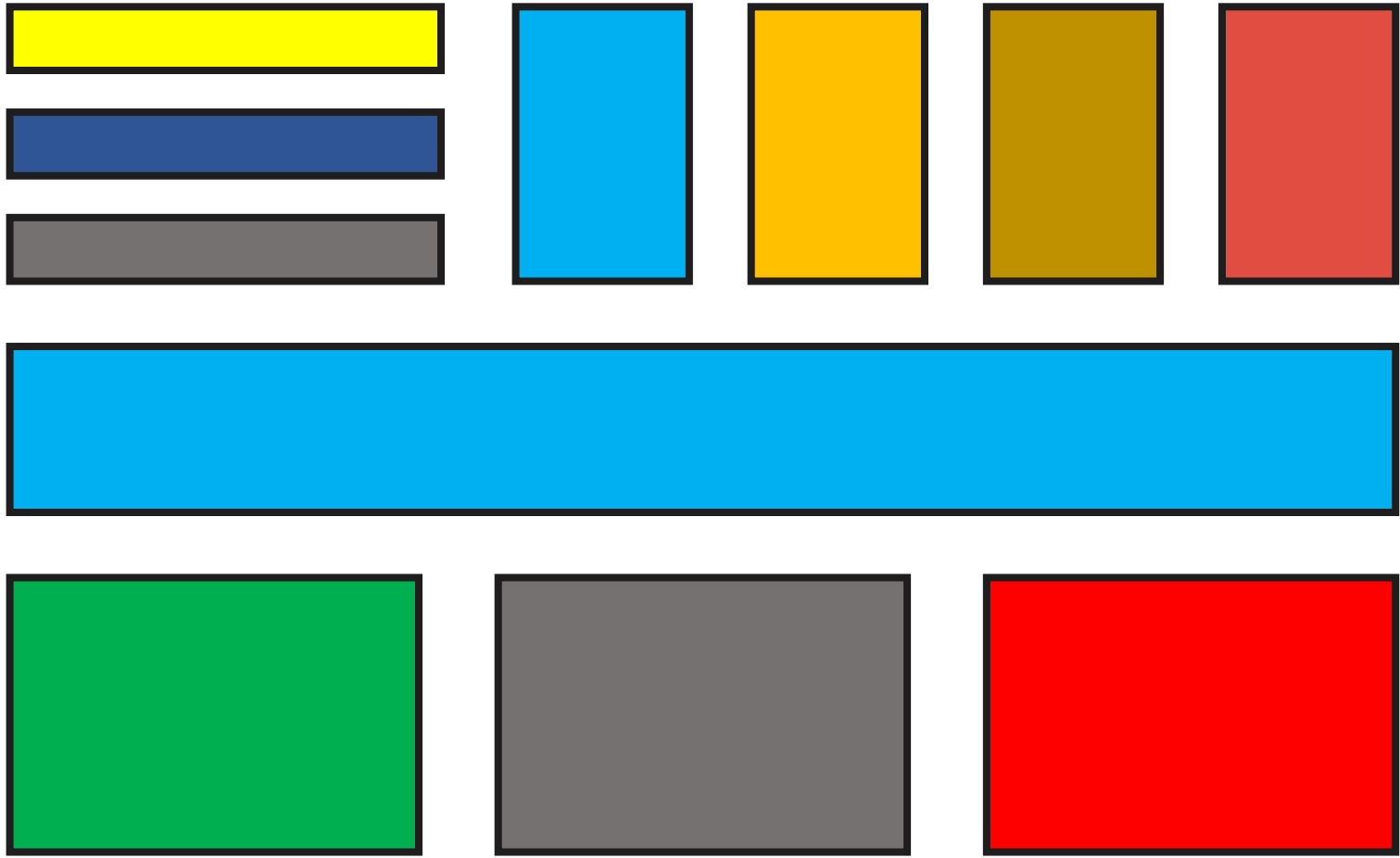




monolith









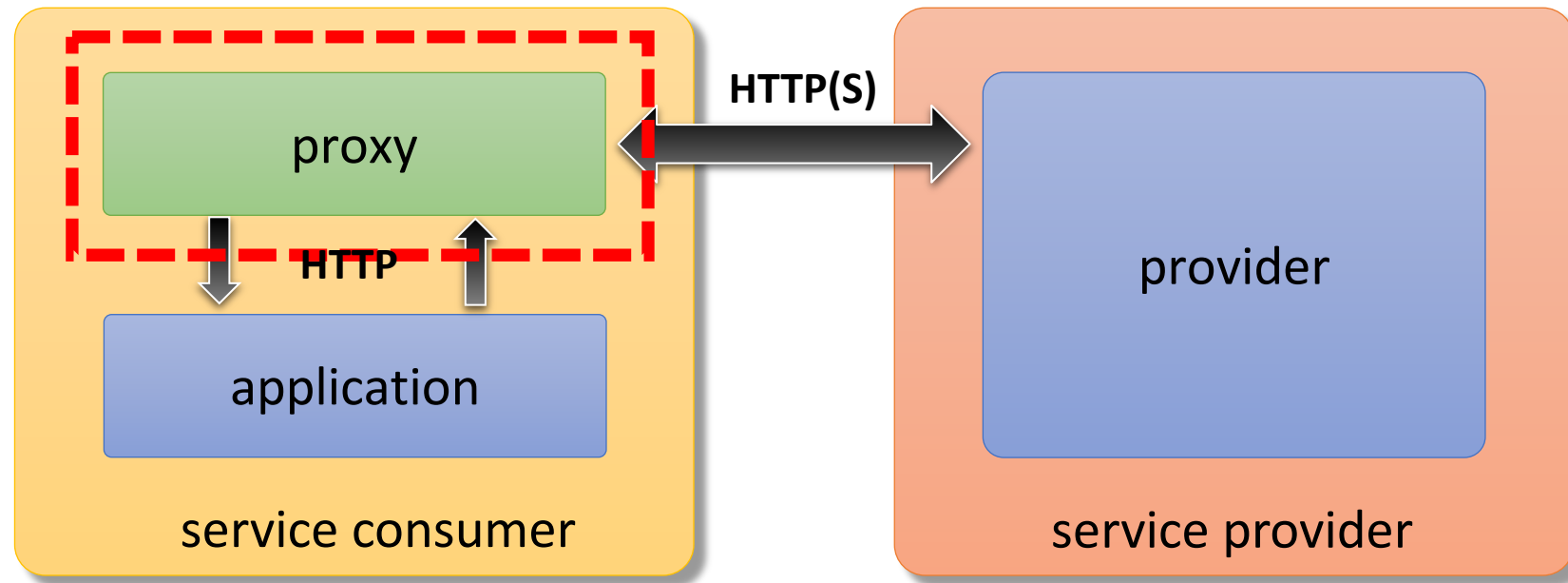
# Consistency & Visibility in communications



# Our setup

.....

data plane





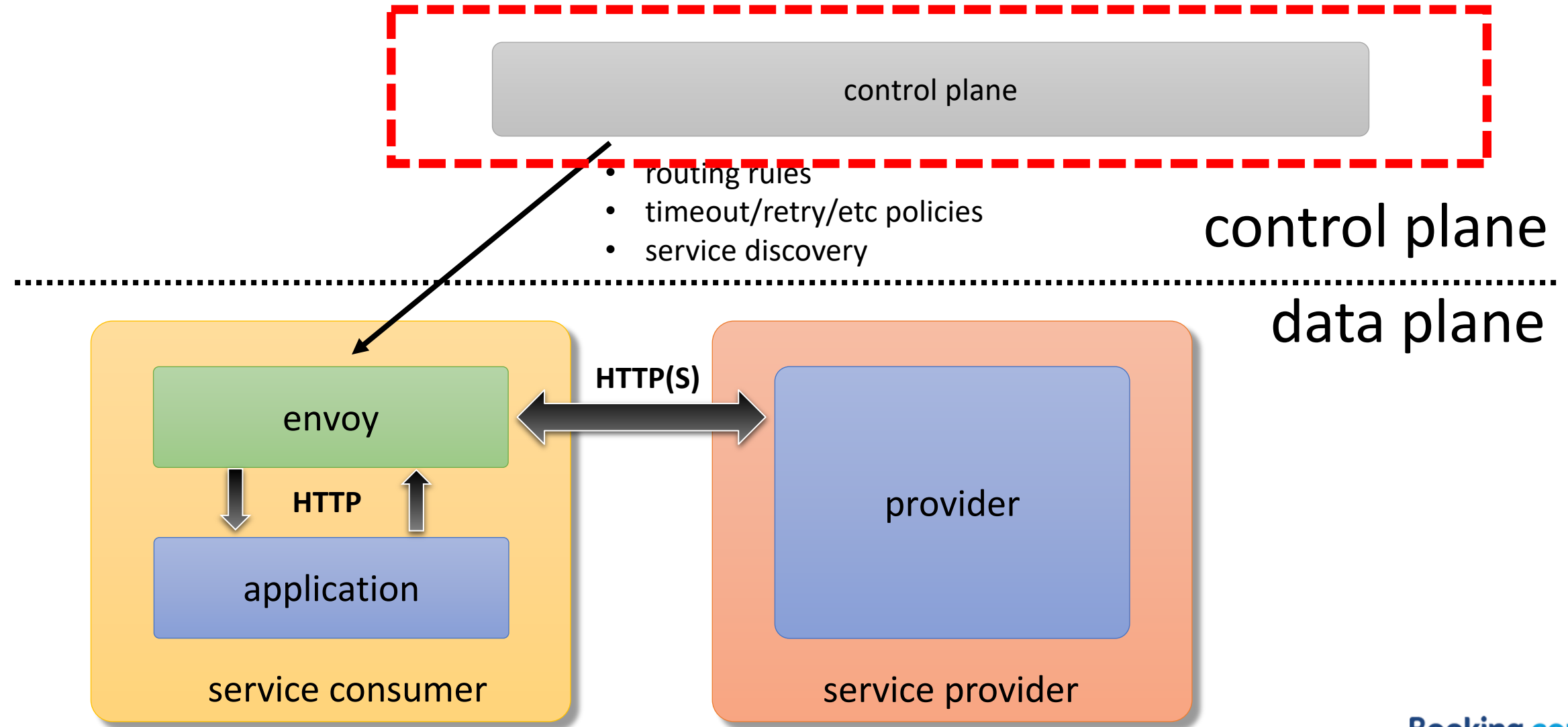
**linkerd**



**envoy**

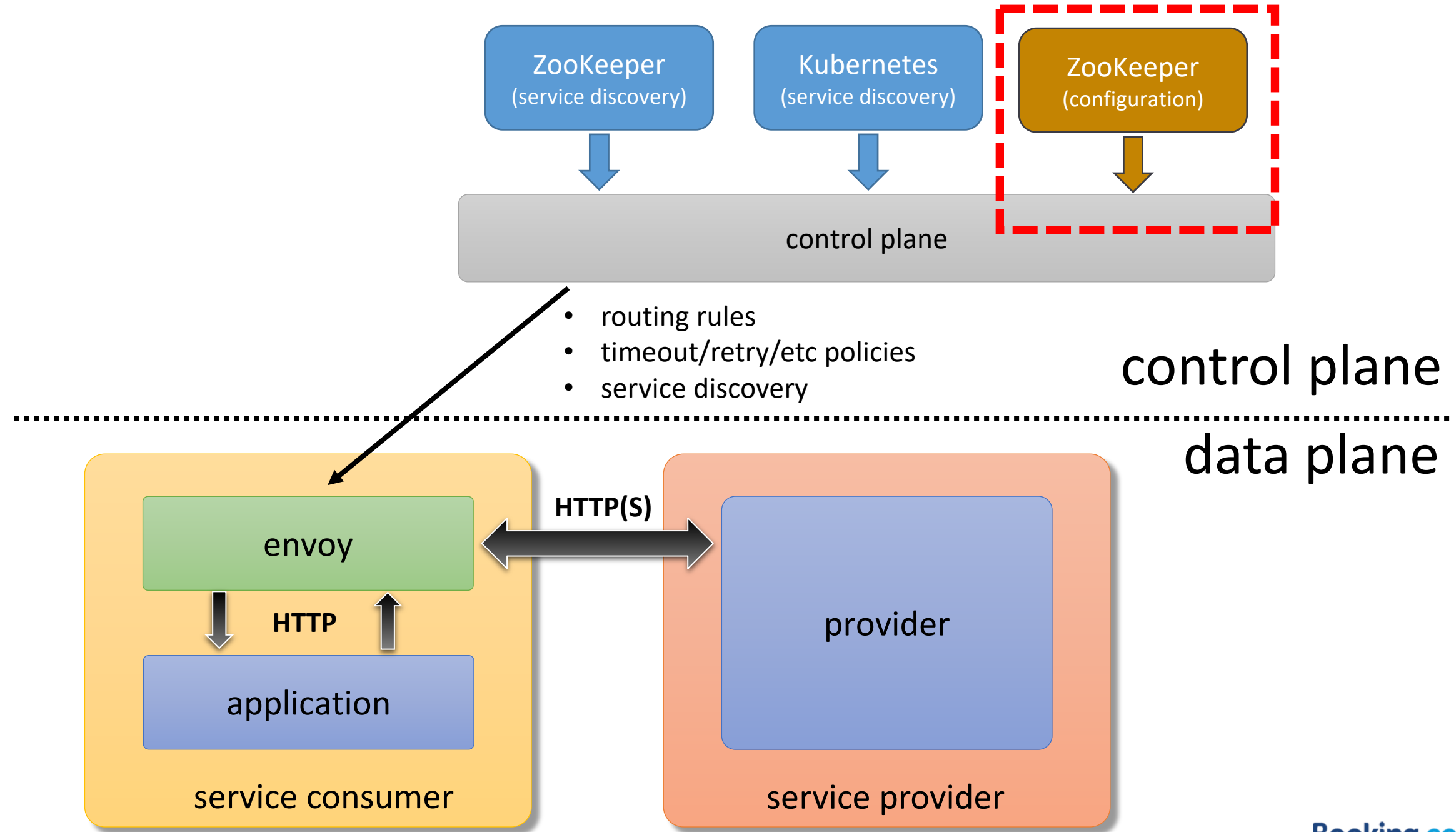
- graceful restart
- TCP proxy





# control plane

- in-house (v1/v2 API)
  - started in August 2017; Istio is around 0.2
  - one complex system at a time
- start with bare-metal support
- minimal abstraction
  - yup, just write Envoy config (partially)
  - fine with Envoy's set of features
- self-service for service owners



# configuration

- Envoy configuration is quite rich
- power vs. usability:
  - [cluster specification](#)
  - [virtual host specification](#)

```
kind: ClusterSpec
metadata:
  annotations:
    watcher_name: zookeeper.watcher
    paths: ["/pools/api/dc"]
spec:
  name: api.prod.cluster
  lb_policy: LEAST_REQUEST
  connect_timeout: 1s
  lb_subset_config:
    subset_selectors:
      - keys: ["DC"]
      - keys: ["IsLocal"]
    fallback_policy: DEFAULT_SUBSET
  default_subset:
    IsLocal: true
```

```
kind: VirtualHostSpec
spec:
  name: api.vhost
  domains: ["api.service"]
  routes:
    - match:
        prefix: /
      route:
        cluster: api.prod.cluster
        timeout: 10s
        retry_policy:
          retry_on: 5xx
          num_retries: 3
          per_try_timeout: 3s
```

**kind: ClusterSpec**

metadata:

annotations:

watcher\_name: zookeeper.watcher

paths: ["/pools/api/dc"]

spec:

name: api.prod.cluster

lb\_policy: LEAST\_REQUEST

connect\_timeout: 1s

lb\_subset\_config:

subset\_selectors:

- keys: ["DC"]

- keys: ["IsLocal"]

fallback\_policy: DEFAULT\_SUBSET

default\_subset:

IsLocal: true

**kind: VirtualHostSpec**

spec:

name: api.vhost

domains: ["api.service"]

routes:

- match:

prefix: /

route:

cluster: api.prod.cluster

timeout: 10s

retry\_policy:

retry\_on: 5xx

num\_retries: 3

per\_try\_timeout: 3s

```
kind: ClusterSpec
metadata:
  annotations:
    watcher_name: zookeeper.watcher
    paths: ["/pools/api/dc"]
spec:
  name: api.prod.cluster
  lb_policy: LEAST_REQUEST
  connect_timeout: 1s
  lb_subset_config:
    subset_selectors:
      - keys: ["DC"]
      - keys: ["IsLocal"]
    fallback_policy: DEFAULT_SUBSET
    default_subset:
      IsLocal: true
```

```
kind: VirtualHostSpec
spec:
  name: api.vhost
  domains: ["api.service"]
  routes:
    - match:
        prefix: /
      route:
        cluster: api.prod.cluster
        timeout: 10s
        retry_policy:
          retry_on: 5xx
          num_retries: 3
          per_try_timeout: 3s
```

```
kind: ClusterSpec
metadata:
  annotations:
    watcher_name: zookeeper.watcher
    paths: ["/pools/api/dc"]
spec:
  name: api.prod.cluster
  lb_policy: LEAST_REQUEST
  connect_timeout: 1s
  lb_subset_config:
    subset_selectors:
      - keys: ["DC"]
      - keys: ["IsLocal"]
    fallback_policy: DEFAULT_SUBSET
  default_subset:
    IsLocal: true
```

```
kind: VirtualHostSpec
spec:
  name: api.vhost
  domains: ["api.service"]
  routes:
    - match:
        prefix: /
      route:
        cluster: api.prod.cluster
        timeout: 10s
        retry_policy:
          retry_on: 5xx
          num_retries: 3
          per_try_timeout: 3s
```



```
kind: ClusterSpec
metadata:
  annotations:
    watcher_name: zookeeper.watcher
    paths: ["/pools/api/dc"]
```

spec:

```
name: api.prod.cluster
lb_policy: LEAST_REQUEST
connect_timeout: 1s
lb_subset_config:
  subset_selectors:
  - keys: ["DC"]
  - keys: ["IsLocal"]
fallback_policy: DEFAULT_SUBSET
default_subset:
  IsLocal: true
```

envoy cluster spec

```
kind: VirtualHostSpec
spec:
```

```
name: api.vhost
domains: ["api.service"]
routes:
- match:
  prefix: /
  route:
    cluster: api.prod.cluster
    timeout: 10s
    retry_policy:
      retry_on: 5xx
      num_retries: 3
      per_try_timeout: 3s
```

envoy virtual host spec

```
kind: ClusterSpec
metadata:
  annotations:
    watcher_name: zookeeper.watcher
    paths: ["/pools/api/dc"]
```

```
spec:
```

```
name: api.pr
lb_policy: I
connect_time
lb_subset_co
subset_sel
- keys: ['
- keys: ['
fallback_po
default_subset:
  IsLocal: true
```

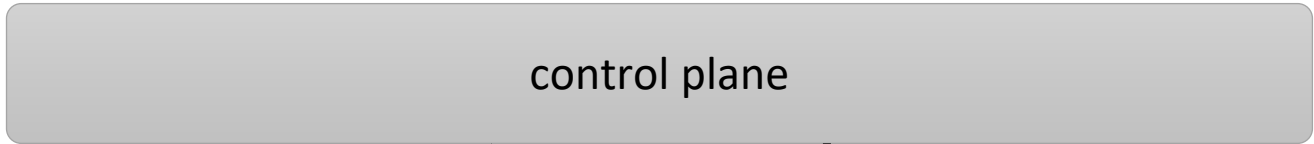
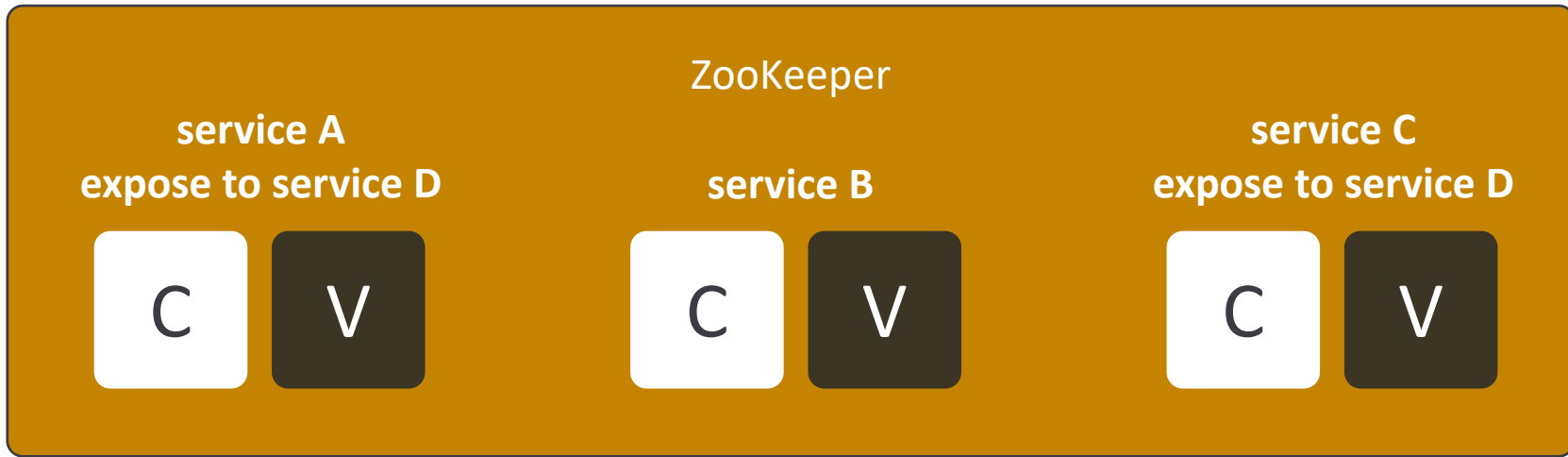
envoy cluster spec

```
kind: VirtualHostSpec
spec:
```

```
name: api.vhost
domains: ["api.service"]
routes:
```

bootstrap wizard

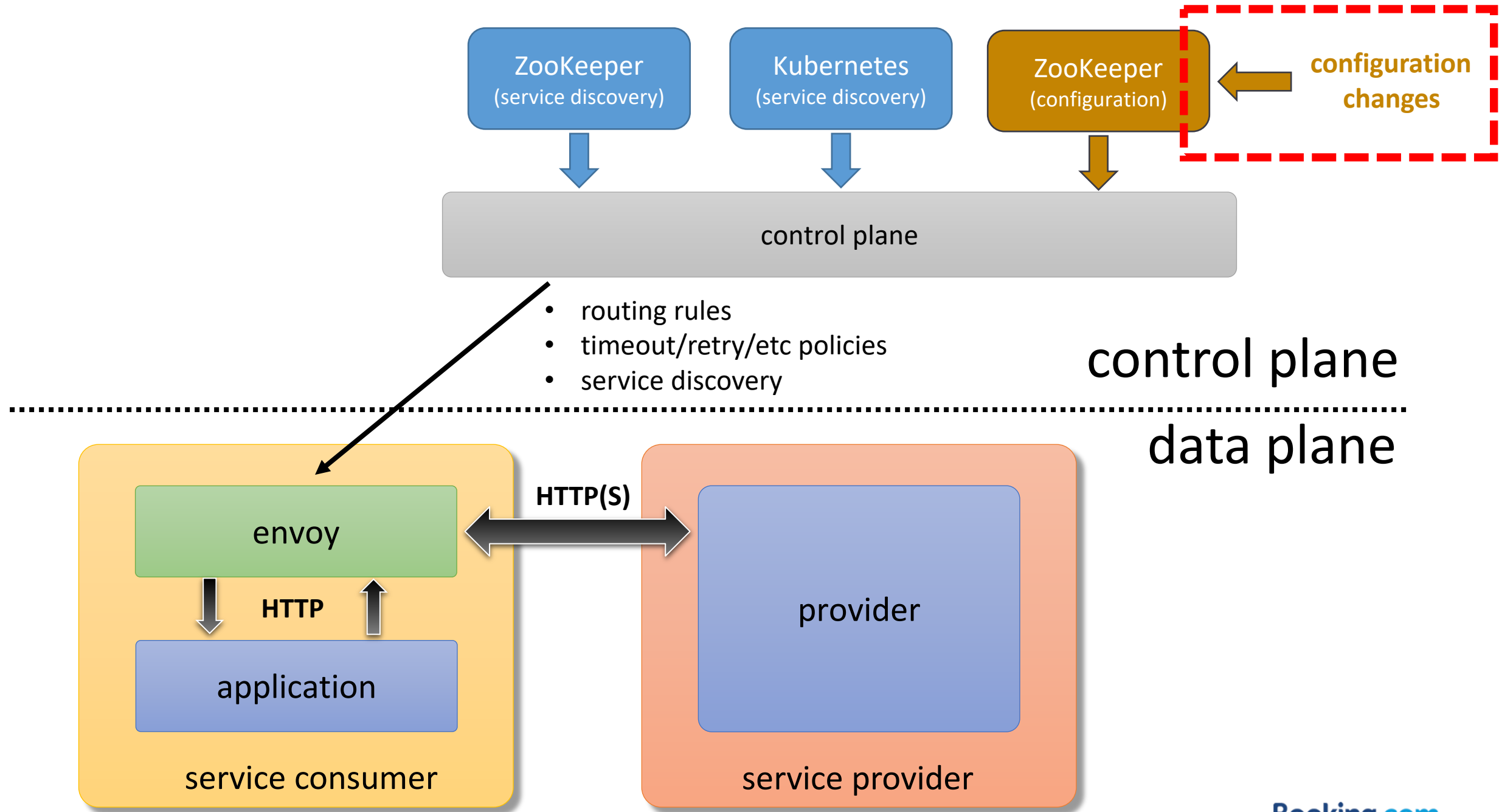
envoy virtual host spec



I'm service D

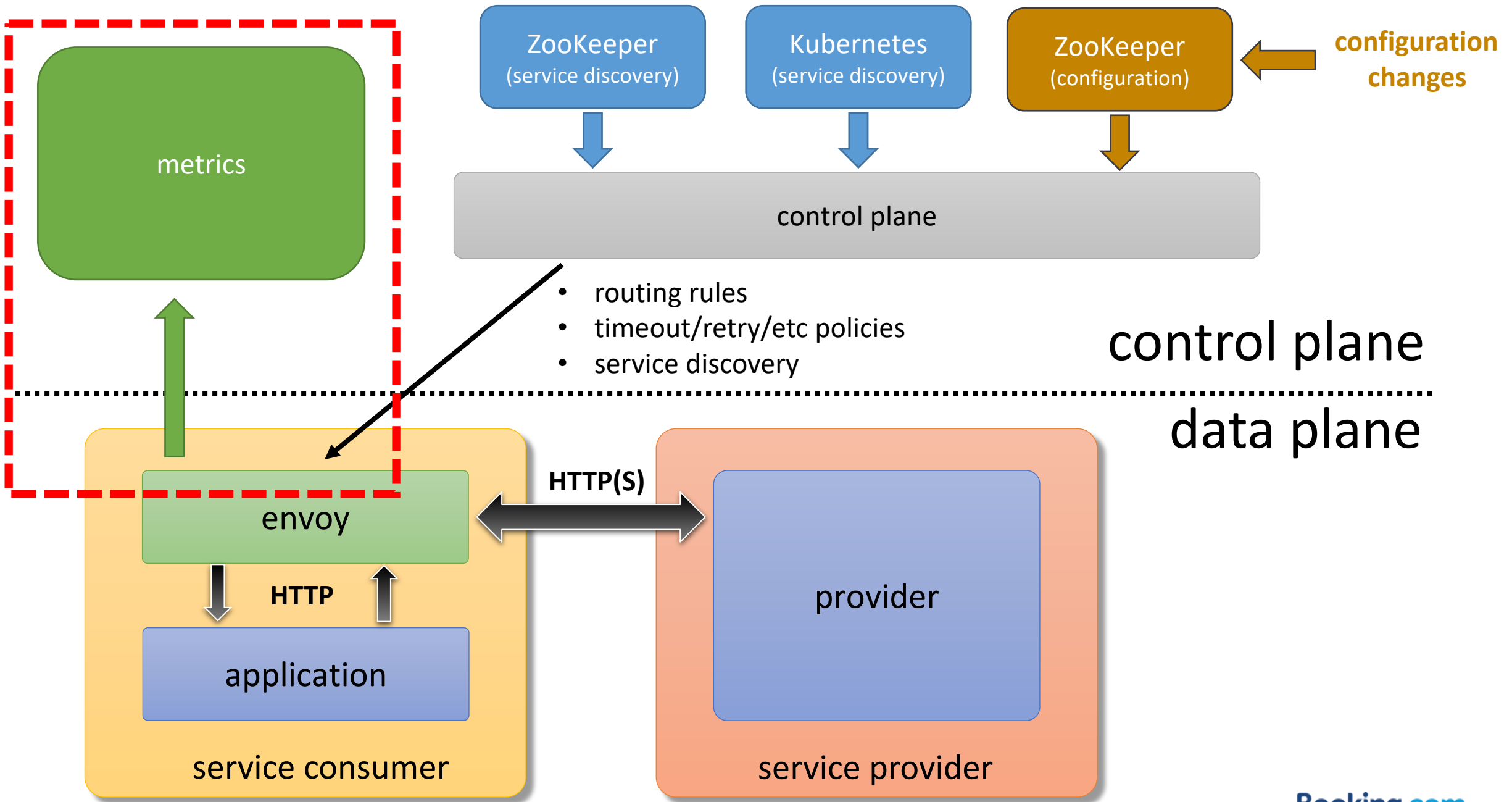


here is service A  
and service C specs



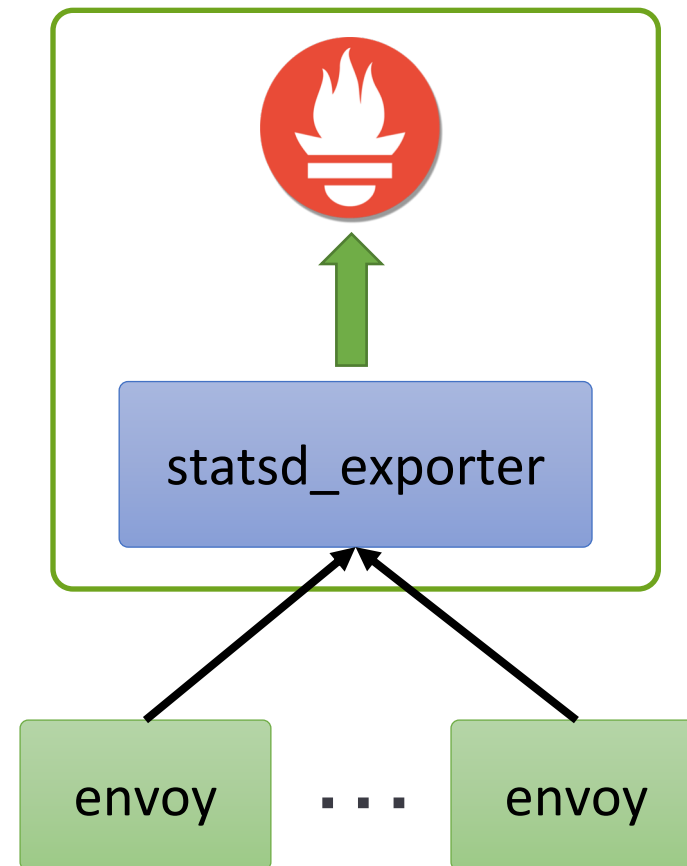
# Deploying changes

- integration with existent infrastructure
  - [git-deploy](#)
- vanguard (canary) deployment
- syntax and semantic validation
- fast rollback

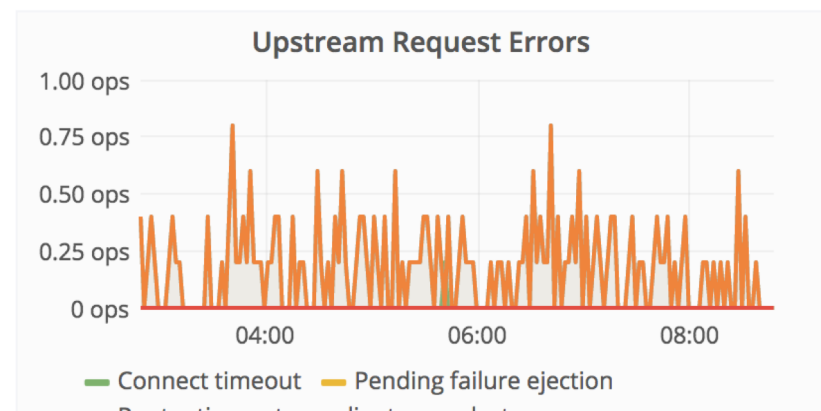
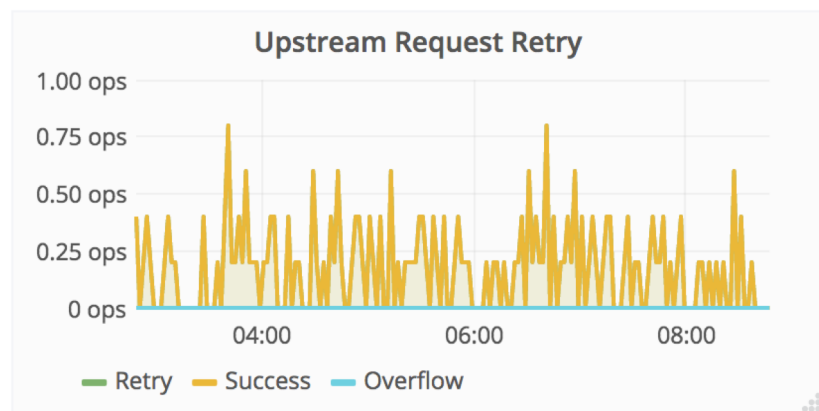
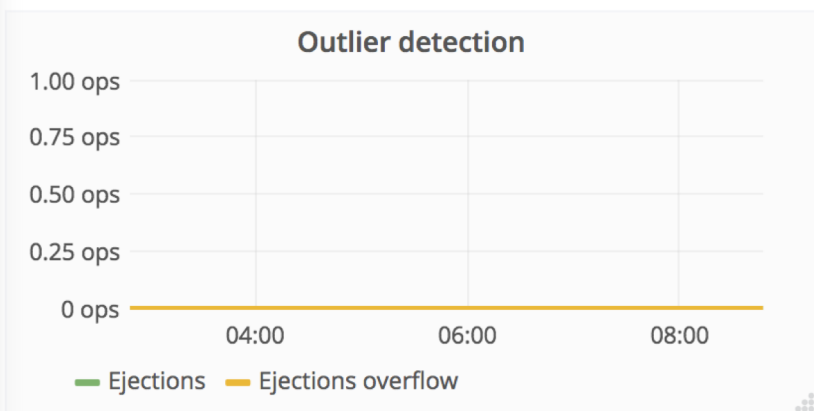
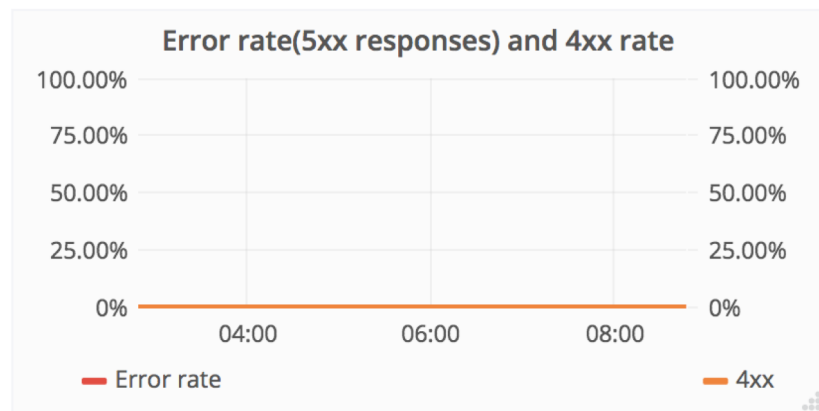
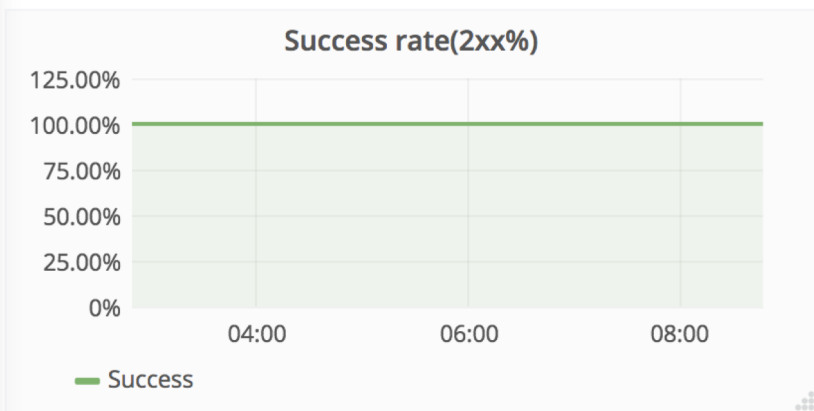
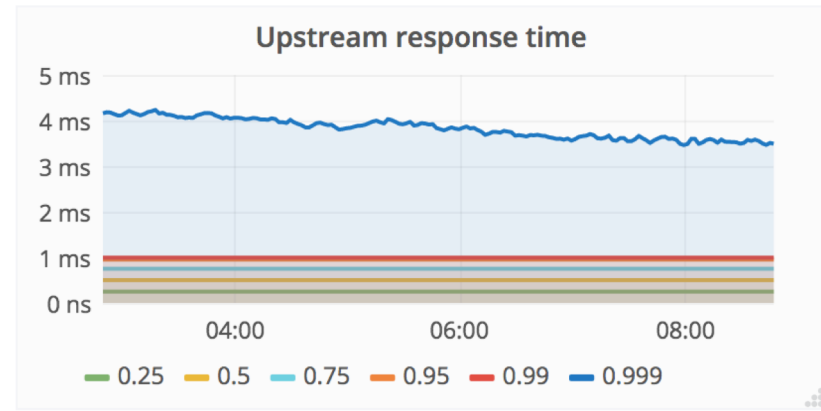
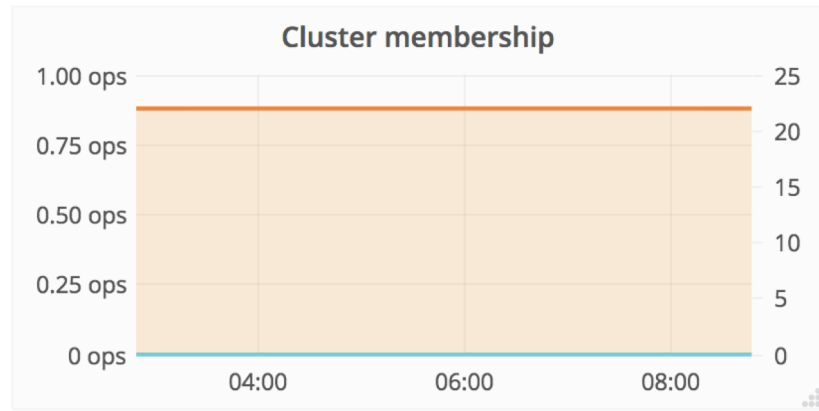
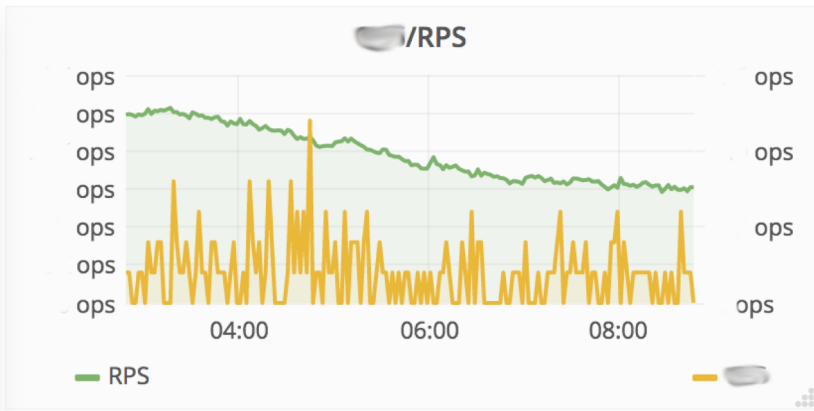


# Monitoring

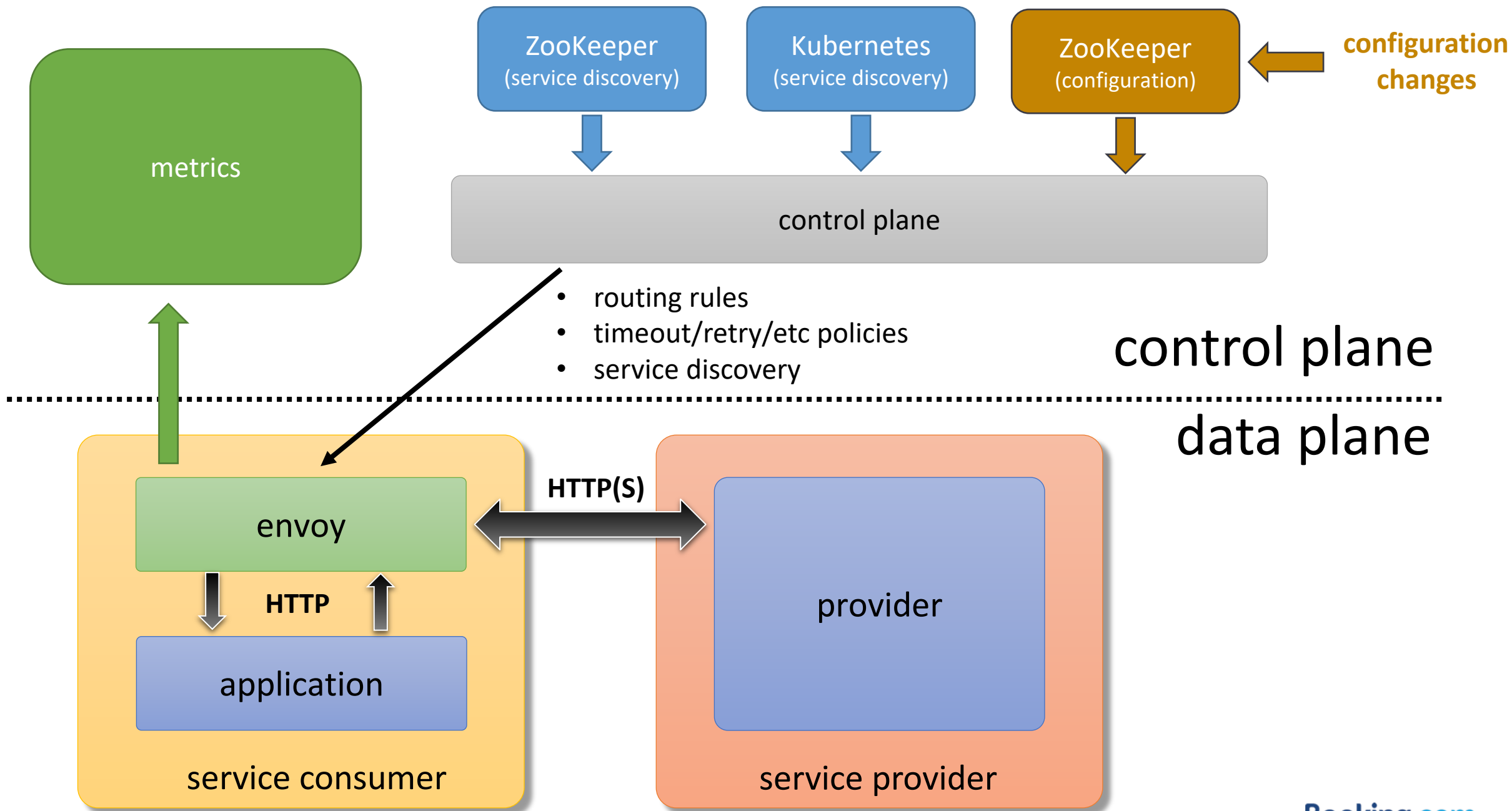
- approach #1 – graphite
  - excellent in-house facilities
  - aggregations is too heavy at scale
- approach #2 – prometheus
  - statsd support
  - statsd\_exporter
  - still iterating
- standard dashboard



▼ Egress to /prod.cluster







# Production

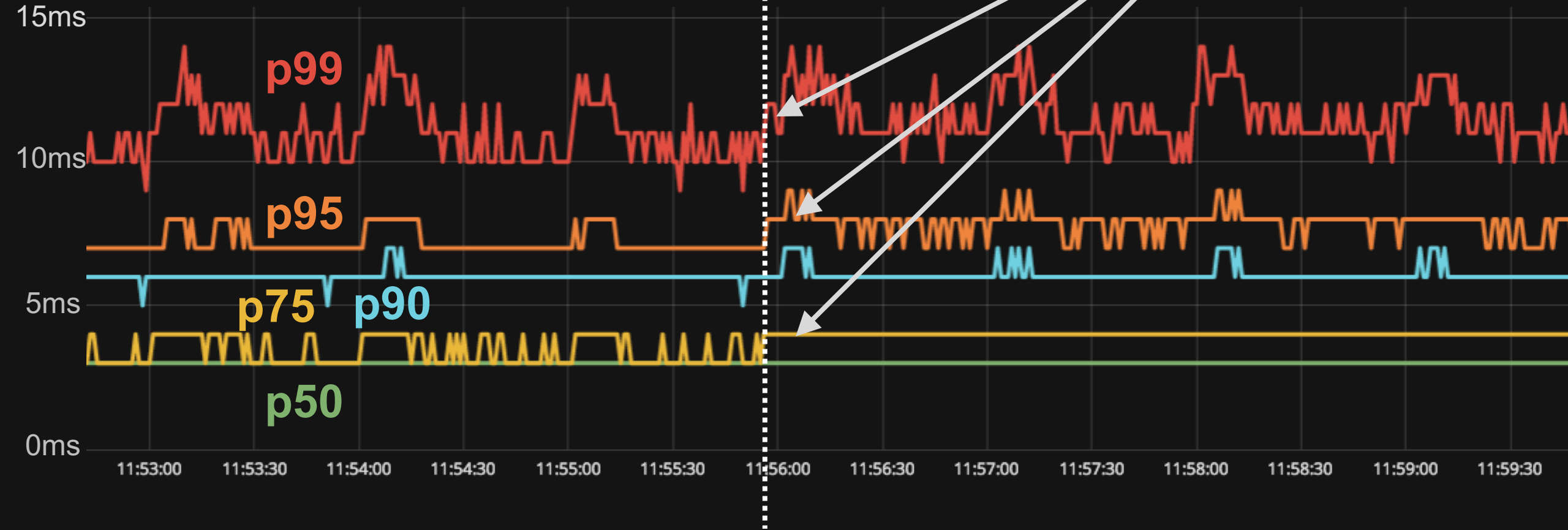
# Some numbers

- ~ 6 months in production
- ~ 30 projects
- ~ 10K servers
- hundreds of thousands RPS
- overheads...

# HTTP

perceived latency

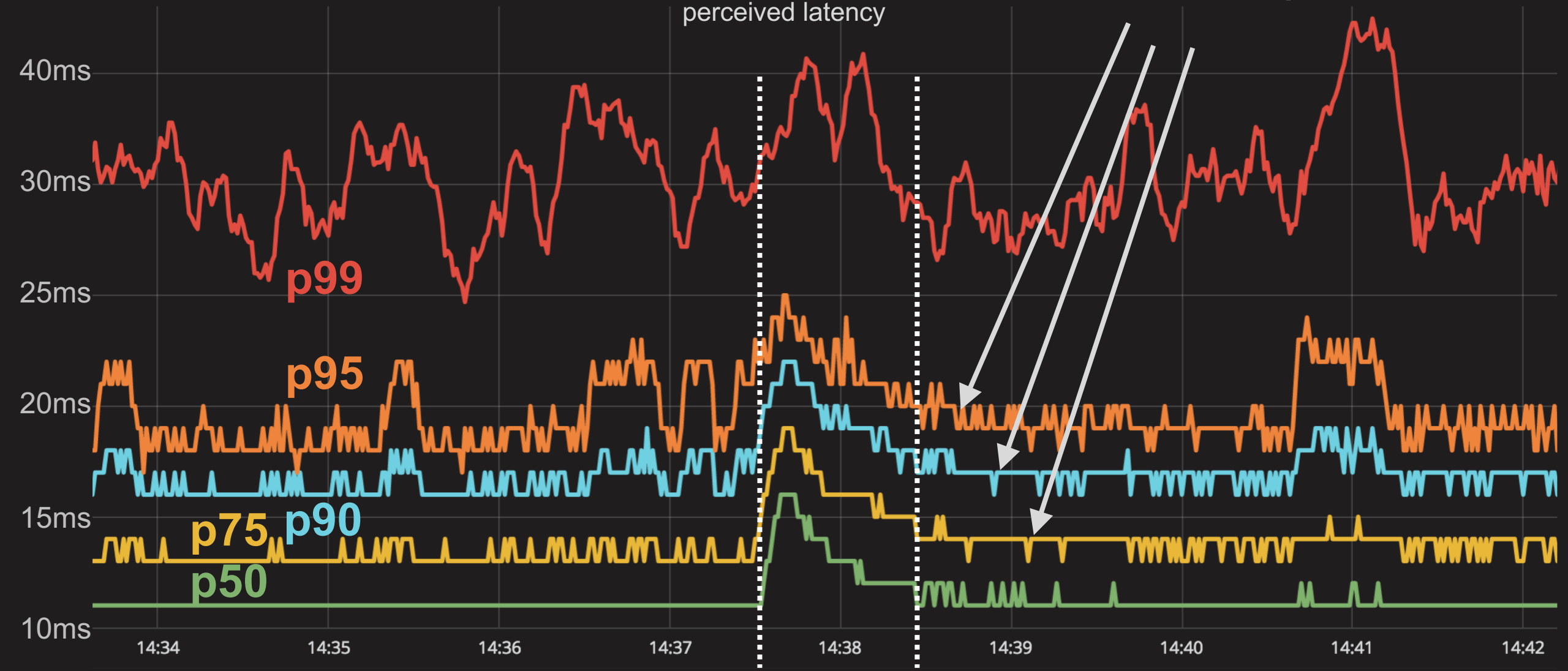
+ 1ms



# HTTPS

perceived latency

## + 1ms



# Some findings

- graceful restart works, but not always!
  - (major) Envoy upgrades may not be graceful
- "x-envoy-retry-on = connect-failure" – too few  
"x-envoy-retry-on = 5xx" – too much
  - [gateway-error](#) – HTTP 502, 503, 504
- absence of TCP keepalive => stale connections
  - [coming in 1.7](#)
  - [idle timeout](#) in TCP proxy in 1.6
- [cluster name \(1.5\) -> cluster names \(1.6\)](#)

# Conclusion

**service mesh\* is  
a building block on the path to SOA  
technically & organizationally**

\* the way I understand it



**Andrei Vereha**

**Envoy production stories**

**Friday, May 4th 14:00**

**Envoy Deep Dive**



**Thank you!**

**Ivan Kruglov**  
**ivan.kruglov@booking.com**

**Booking.com**