



KubeCon



CloudNativeCon

Europe 2018

# Kubervisor: Pod Anomaly Detection

David Benque, Amadeus S.A.S

Cedric Lamoriniere, Amadeus S.A.S



# Who are we?



KubeCon



CloudNativeCon

Europe 2018

## David Benque

@BenqueDavid

Software Engineer at Amadeus

---

CNCF Meetup Organiser

Distributed systems

PaaS

Automation

Gopher

## Cedric Lamoriniere

@cedriclam

Software Engineer at Amadeus

---

CNCF Meetup Organiser

Distributed systems

Gopher

# Amadeus

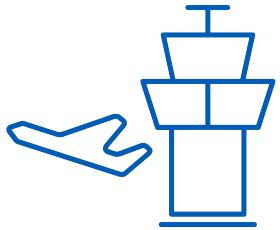


KubeCon

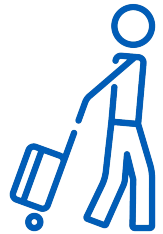


CloudNativeCon

Europe 2018



~600 million total bookings processed in 2016



1.3 billion passengers boarded in 2016



~450 000 queries per second (600 000 at peak)

# Agenda



KubeCon



CloudNativeCon

Europe 2018

## 1 Distributed Systems and Stability

---

### 2 Solutions to increase reliability

- Kubernetes integrated solution
  - Addons solution
- 

### 3 Kubervisor: a pod anomaly detection solution

- Architecture
- Demo



KubeCon



CloudNativeCon

Europe 2018

1

# Distributed Systems and Stability



# Distributed systems and stability



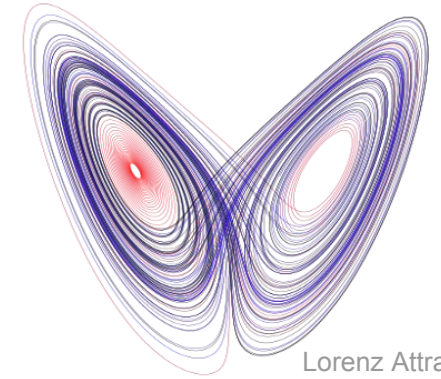
KubeCon



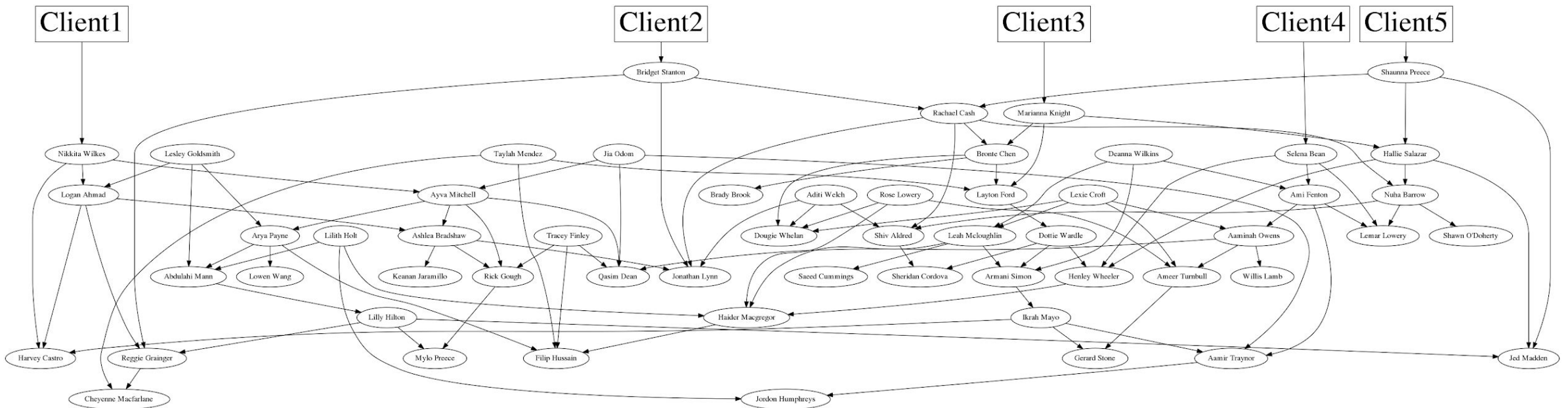
CloudNativeCon

Europe 2018

## Butterfly effect in distributed systems



Lorenz Attractor



# Distributed systems and stability



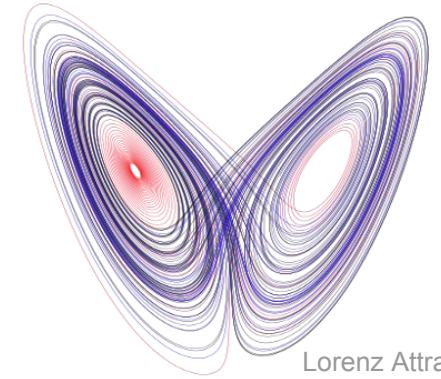
KubeCon



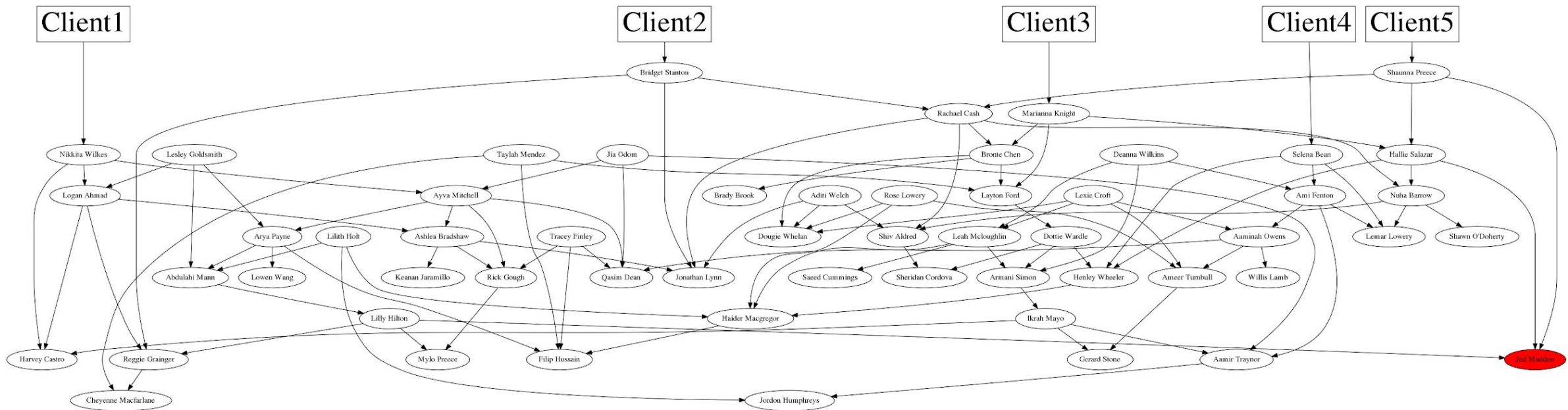
CloudNativeCon

Europe 2018

## Butterfly effect in distributed systems



Lorenz Attractor



# Distributed systems and stability



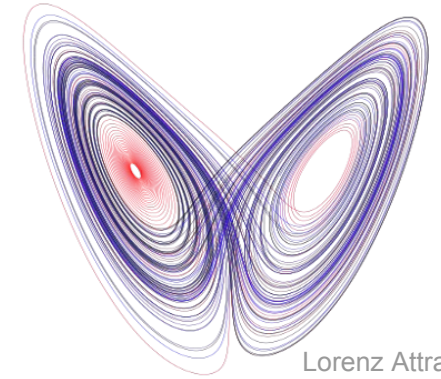
KubeCon



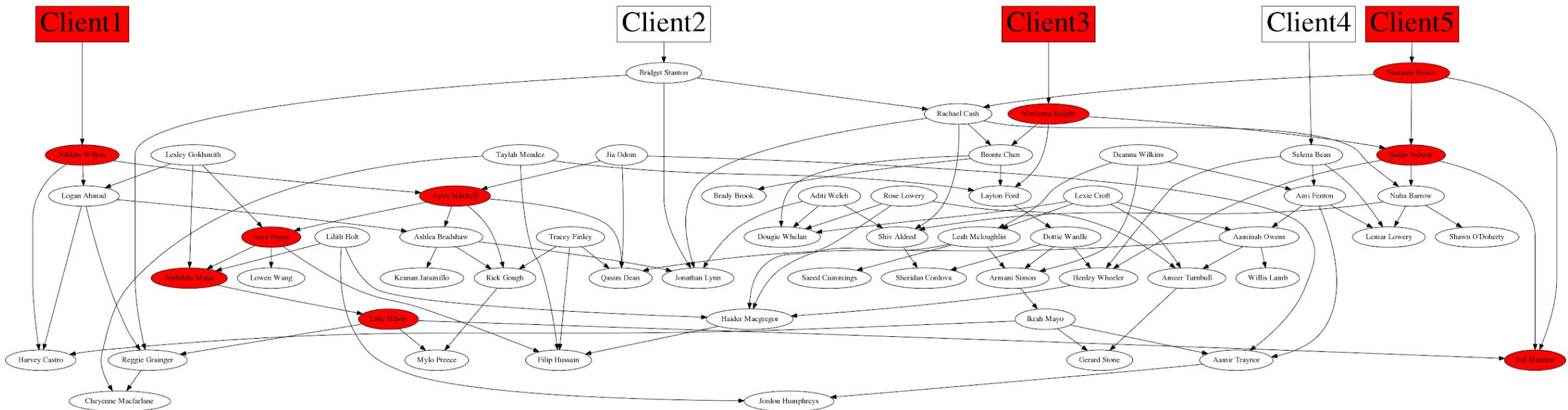
CloudNativeCon

Europe 2018

## Butterfly effect in distributed systems



Lorenz Attractor





# Distributed systems and stability



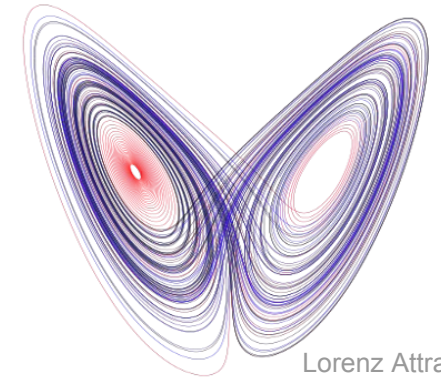
KubeCon



CloudNativeCon

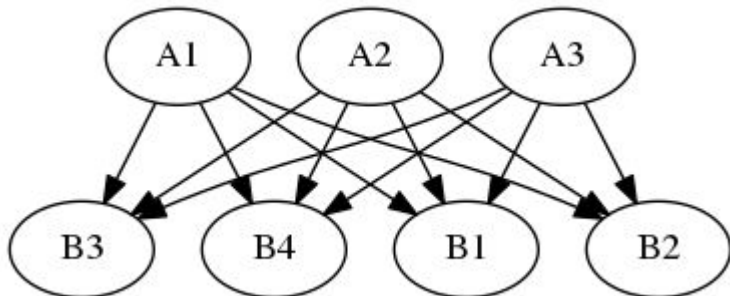
Europe 2018

Load Balancing and multiple requests to a failing dependency could even make things worse

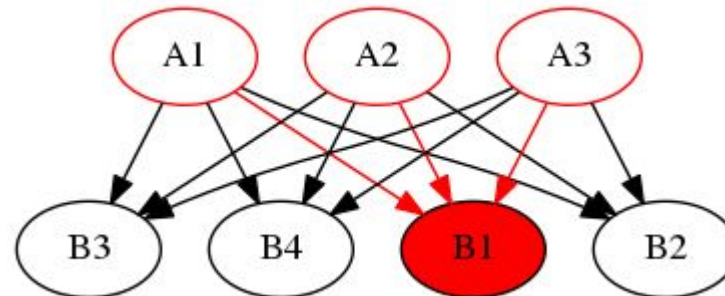


Lorenz Attractor

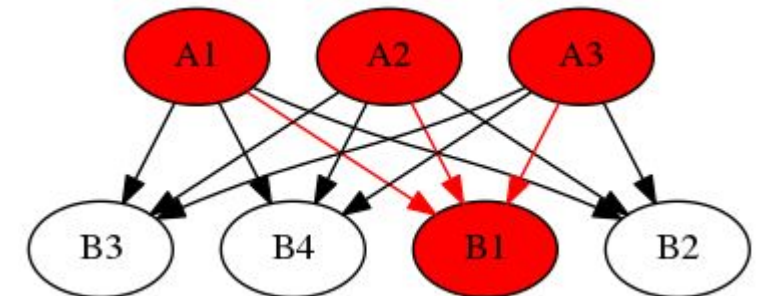
→ All Fine



→ B1 Instance fails  
→ A<sub>i</sub> instances failure rate 25%



→ B1 Instance fails  
→ A<sub>i</sub> makes parallel calls to B<sub>j</sub>  
→ A<sub>i</sub> instances failure rate >25%



# Distributed systems and stability



KubeCon

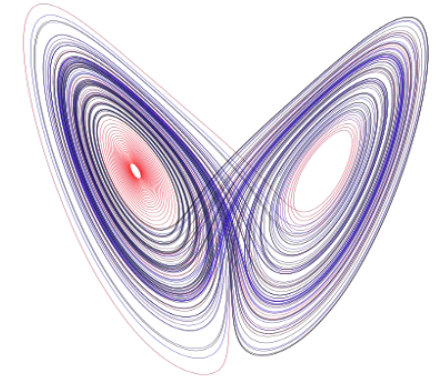


CloudNativeCon

Europe 2018

“A distributed system is one that prevents you from working because of the failure of a machine that you had never heard of.”

Leslie Lamport



Attractor

Edward Lorenz

“We live in a rainbow of chaos.”

Paul Cezanne



**KubeCon**



**CloudNativeCon**

Europe 2018

2

Solutions to increase reliability



# Solutions to increase reliability

## Proximity-based Load Balancing



KubeCon



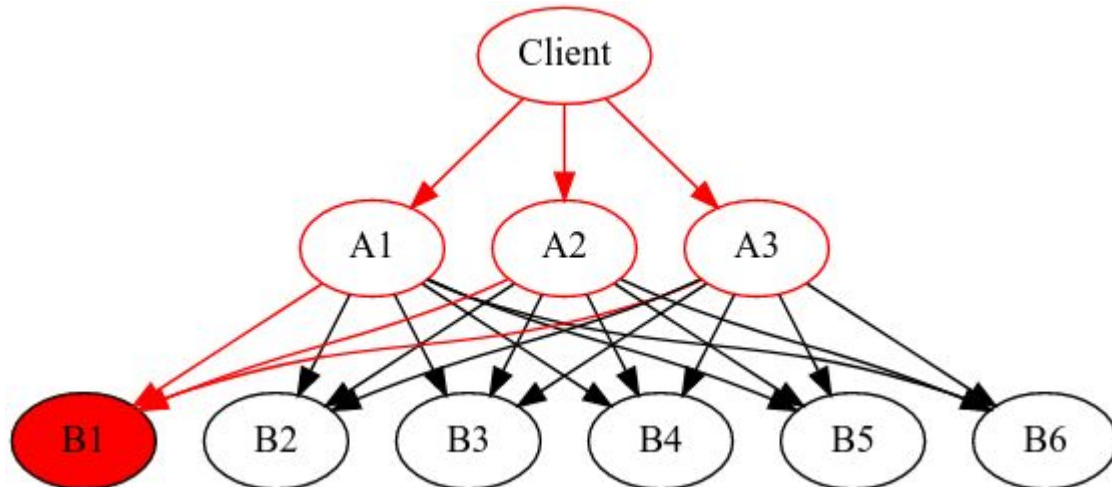
CloudNativeCon

Europe 2018

Assumption: loadbalancing using round robin

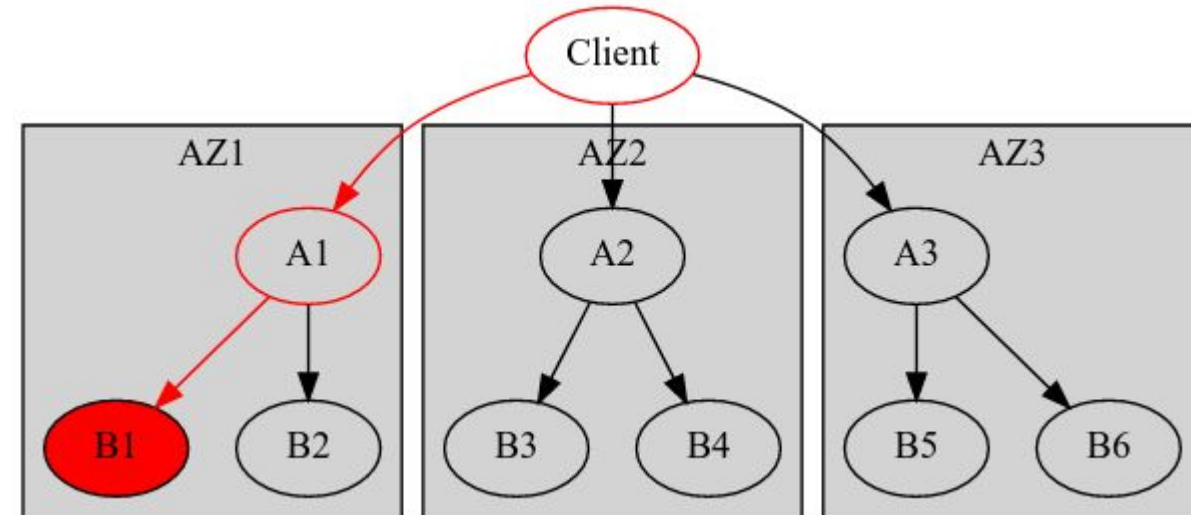
Client->A->B: failure rate  $\frac{1}{6}$

Client->A-(x6)->B: failure rate 1



Client->A->B: failure rate  $\frac{1}{6}$

Client->A-(x6)->B: failure rate  $\frac{1}{3}$



- Ease root cause analysis
- Safeguards overall success rate by constraining fan out

# Solutions to increase reliability

## Proximity-based Load Balancing



KubeCon



CloudNativeCon

Europe 2018

kube-proxy loadbalancing: Use local pods only #7433

 Open evpapp opened this issue on Apr 28, 2015 · 32 comments

Topology aware service routing #58732

 Open m1093782566 wants to merge 4 commits into `kubernetes:master` from `m1093782566:service-topology`

proposal: topology aware routing of services #1551

 Open m1093782566 wants to merge 6 commits into `kubernetes:master` from `m1093782566:service-topologykey`

kubernetes/kubernetes

kubernetes/community



pilot-agent proxy (envoy) is Availability Zone aware:

`--availabilityZone <string>` : The availability zone where this Envoy instance is running. When running Envoy as a sidecar in Kubernetes, this flag must be one of the availability zones assigned to a node using `failure-domain.beta.kubernetes.io/zone` annotation.

# Solutions to increase reliability

## Container Termination



KubeCon



CloudNativeCon

Europe 2018

“ A Container can exceed its memory request if the Node has memory available. But a Container is **not allowed to use more than its memory limit**. If a Container allocates more memory than its limit, the Container becomes a candidate for termination.”

```
apiVersion: v1
kind: Pod
metadata:
  name: memory-demo-2
  namespace: mem-example
spec:
  containers:
  - name: memory-demo-2-ctr
    image: memtest
    resources:
      requests:
        memory: "50Mi"
      limits:
        memory: "100Mi"
```



<https://kubernetes.io/docs/concepts/configuration/manage-compute-resources-container/>

# Solutions to increase reliability

## Probes

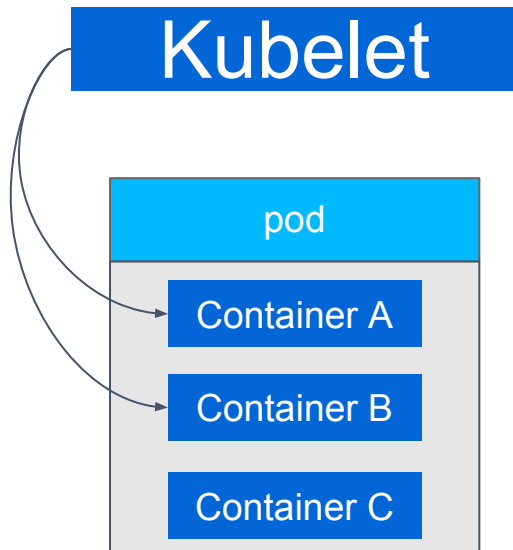


KubeCon



CloudNativeCon

Europe 2018



### Probe Implementation:

- tcpSocket
- httpGet
- exec

```
apiVersion: v1
kind: Pod
...
spec:
  containers:
  - name: A
    ...
    readinessProbe:
      tcpSocket:
        port: 8080
      initialDelaySeconds: 5
    livenessProbe:
      httpGet:
        path: /healthz
        port: 8080
        periodSeconds: 20
  - name: B
    ...
    readinessProbe:
      exec:
        command:
        - myscript
    ...
```

**Liveness** → kill container

**Readiness** → mark endpoint as not ready

// Number of seconds after the container has started before liveness probes are initiated.

`InitialDelaySeconds int32`

// Number of seconds after which the probe times out. (default 1)

`TimeoutSeconds int32`

// How often (in seconds) to perform the probe. (default 10)

`PeriodSeconds int32`

// Minimum consecutive successes for the probe to be considered successful after having failed.(default 1)

`SuccessThreshold int32`

// Minimum consecutive failures for the probe to be considered failed after having succeeded. (default 3)

`FailureThreshold int32`



# Solutions to increase reliability

## Probes



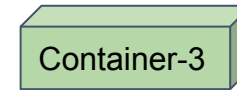
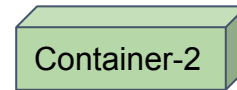
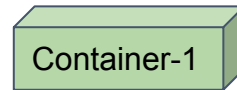
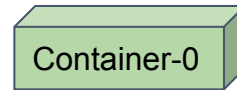
KubeCon



CloudNativeCon

Europe 2018

- Keep it simple.



Complexity → Bugs



# Solutions to increase reliability

## Probes



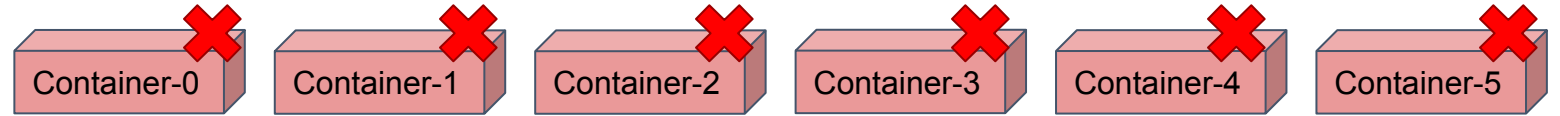
KubeCon



CloudNativeCon

Europe 2018

- Keep it simple.



Complexity → Bugs → Containers collective suicide  
→ Service Outage

# Solutions to increase reliability

## Probes



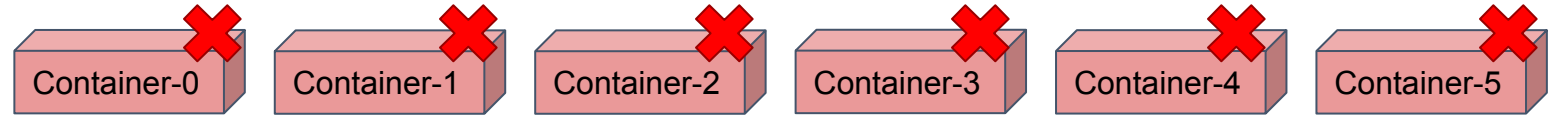
KubeCon



CloudNativeCon

Europe 2018

- Keep it simple.



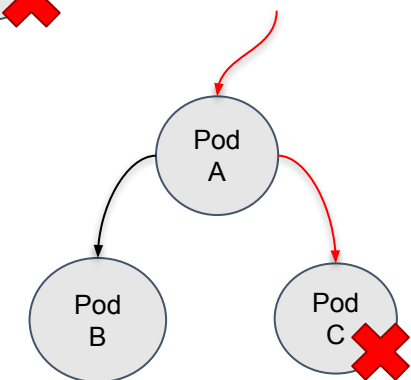
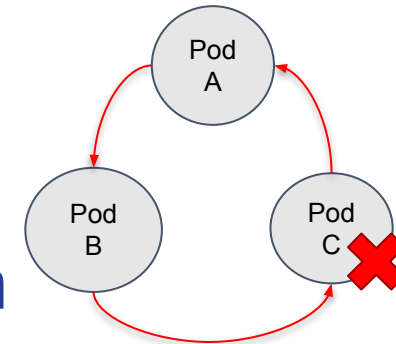
Complexity → Bugs → Containers collective suicide  
→ Service Outage

- Don't check external dependency chain

→ May not be able to restart from scratch

→ Difficult to find problem root cause

→ Failure propagation cutting valid branches

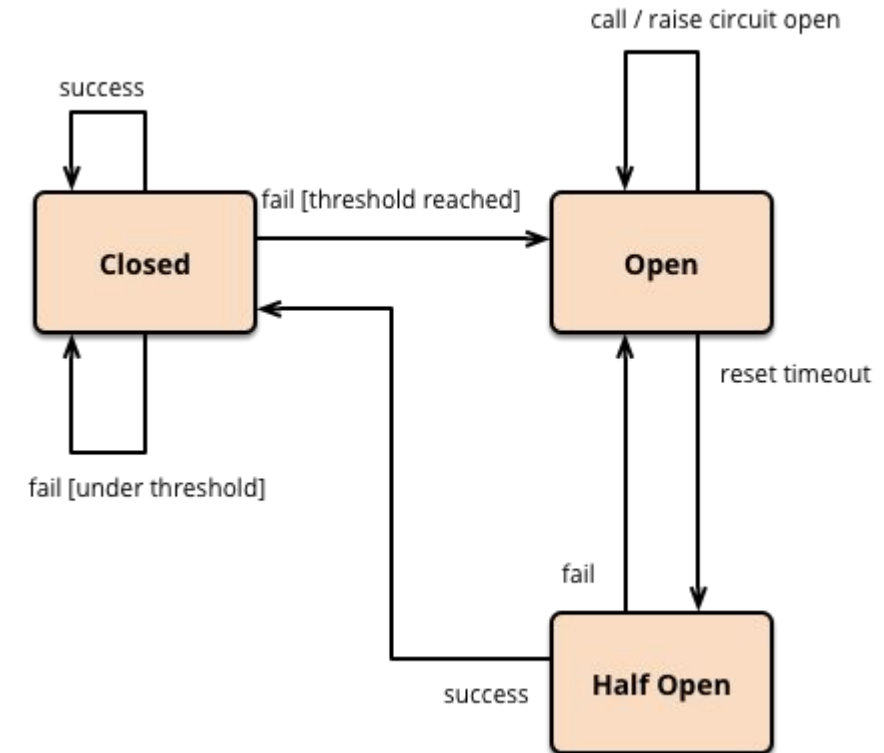
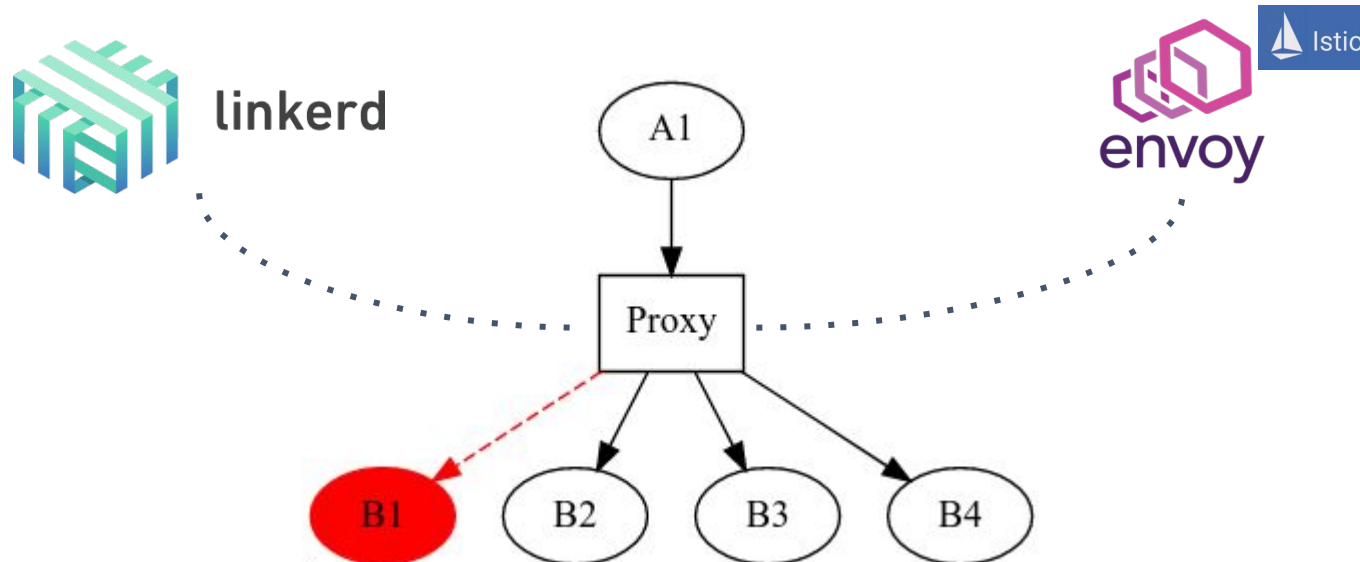


# Solutions to increase reliability

## Circuit Breaker

<https://martinfowler.com/bliki/CircuitBreaker.html>

- The service mesh proxy implements the circuit breaker



# Solutions to increase reliability

## Retries



KubeCon



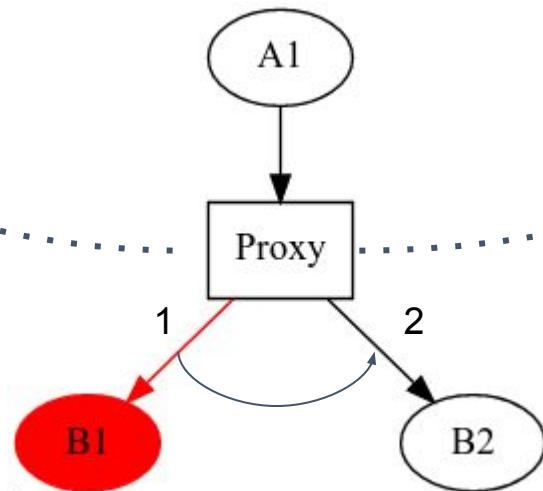
CloudNativeCon

Europe 2018

- The service mesh proxy implements retry policy

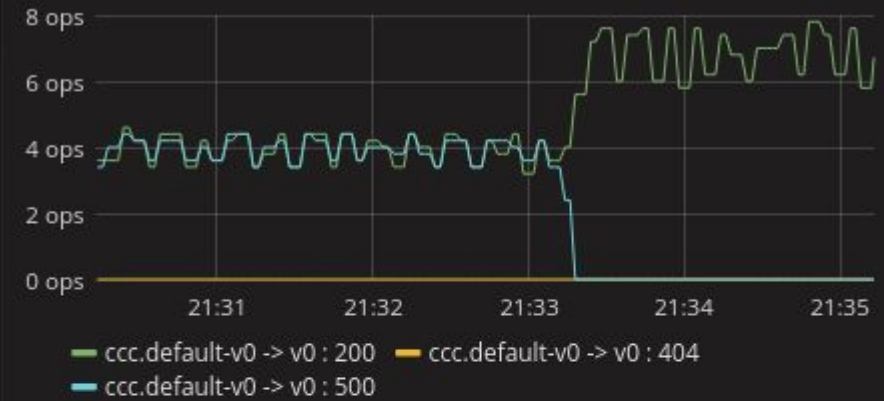


linkerd



aaa.default.svc.cluster.local

Requests by Source, Version, and Response Code



bbb.default.svc.cluster.local

Requests by Source, Version, and Response Code



# Solutions to increase reliability

## Limitations



KubeCon



CloudNativeCon

Europe 2018

- Local decision
  - Probes: in each container
  - Service Mesh: in each proxy
- Based on technical signals only
  - Memory consumption
  - Connection
  - Response time
  - Return Code
  - Circuit Breaker working with HTTP only



KubeCon



CloudNativeCon

Europe 2018

3

## Kubervisor: a pod anomaly detection solution



# Kubervisor

## A pod anomaly detection solution



KubeCon



CloudNativeCon

Europe 2018

- **Orchestrated decision**
  - decision at service level
  - avoid collective suicide, service outage
- **Decision based on technical or/and business information**
  - based on metrics
  - currently supports Prometheus with PromQL.
- **Kubervisor controller:**
  - divided in 2 components: **Breaker, Activator**
  - different strategies possible for the **Breaker and Activator**
  - configuration in dedicated CRD: **KubervisorService**

# Kubervisor

## A pod anomaly detection solution



KubeCon



CloudNativeCon

Europe 2018

- **KubervisorService CRD**

- **spec**

- **service:** application service to watch

- **breaker:** breaker configuration

- global settings
- different implementations

- **activator:** activator configuration

- **mode:** periodic, retryAndPause, retryAndKill
- **period:** pause duration

```
apiVersion: breaker.kubervisor.io/v1
kind: KubervisorService
metadata:
  name: foo
spec:
  service: foo-svc
  breaker:
    minPodsAvailableCount: 2
    minPodsAvailableRatio: 50
    customService: my-metrics-svc:8080
  activator:
    mode: periodic
    period: 15
```

```
maxRetryCount: 2
```



# Kubervisor workflow

## how to use it

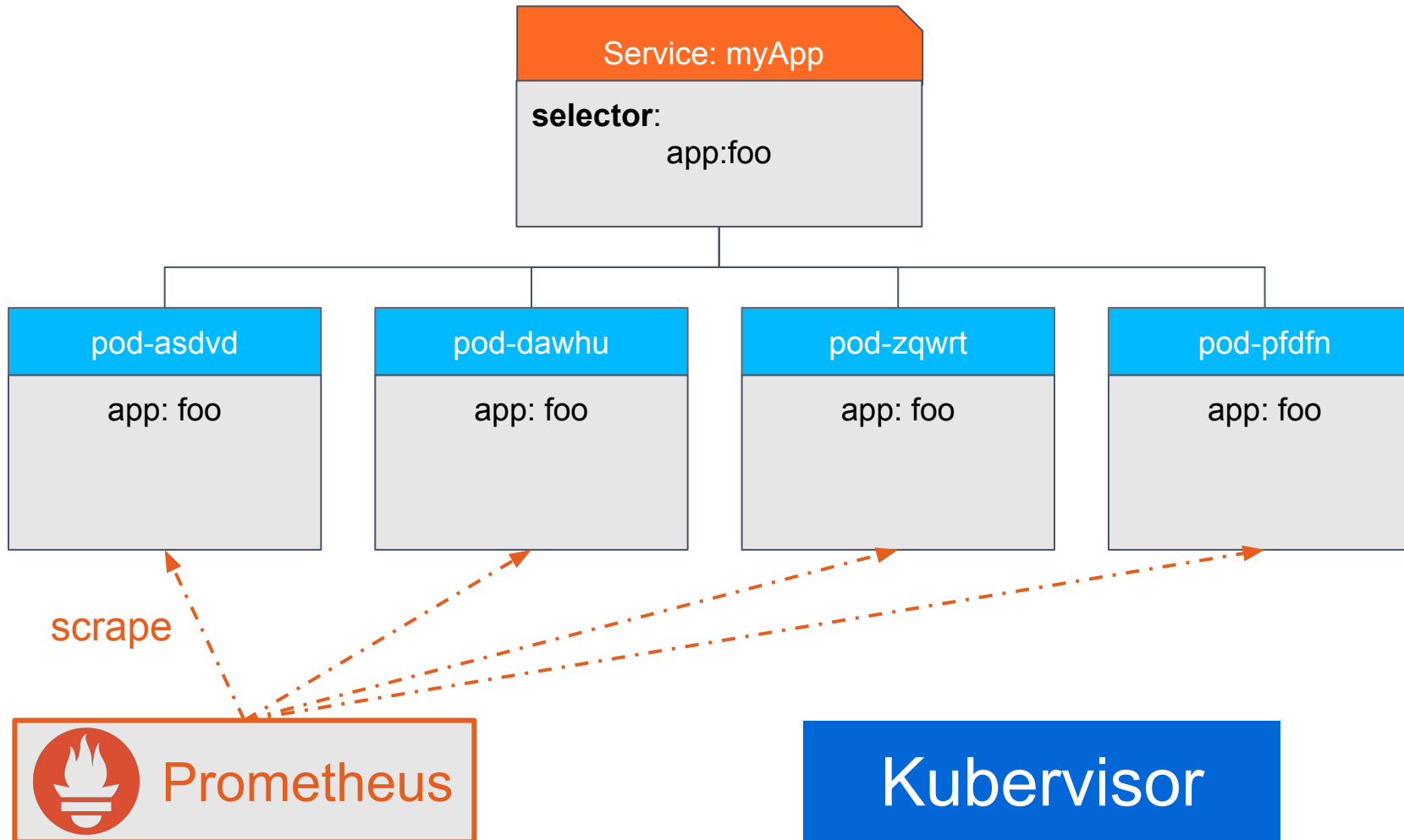


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor workflow

## Initialisation

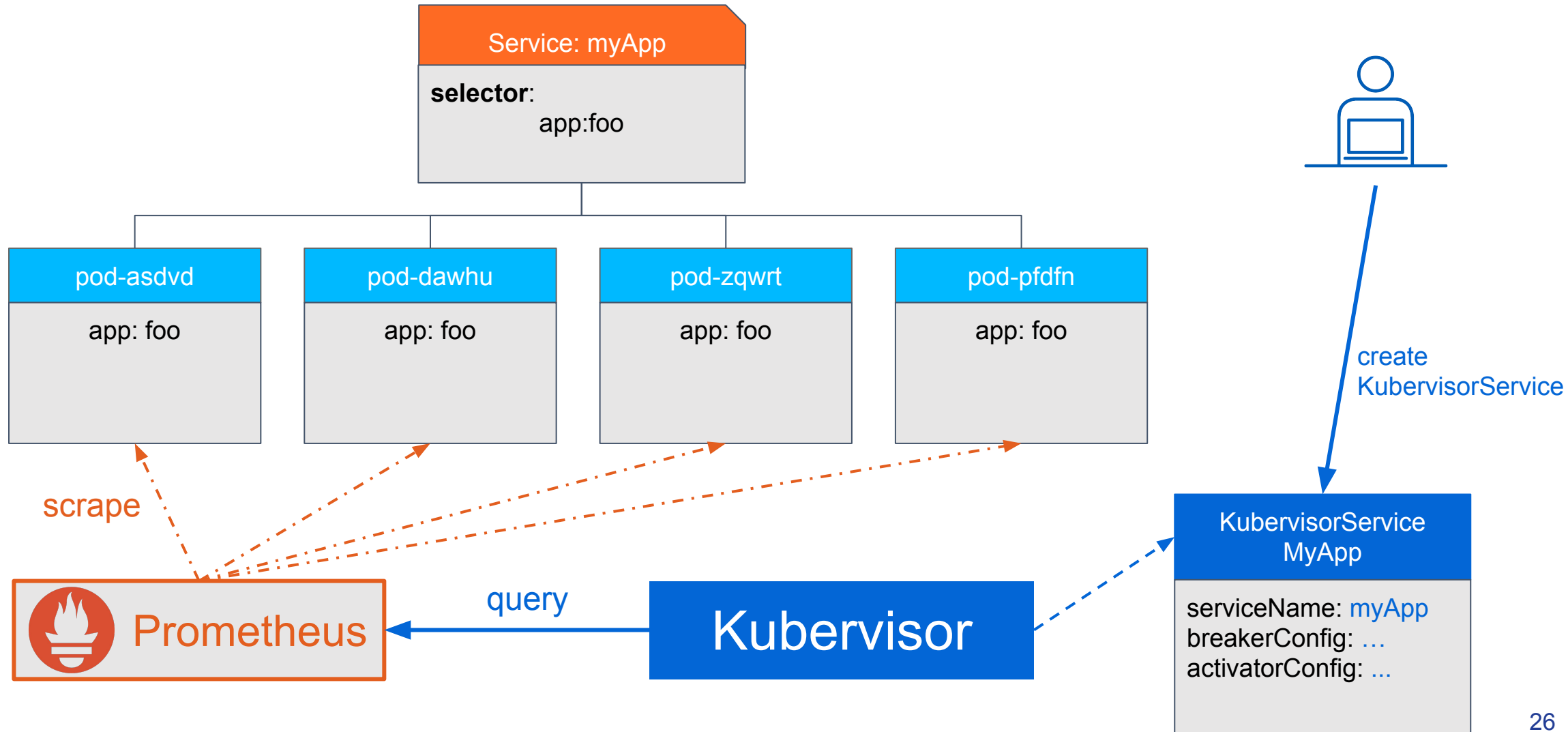


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor workflow

## Initialisation

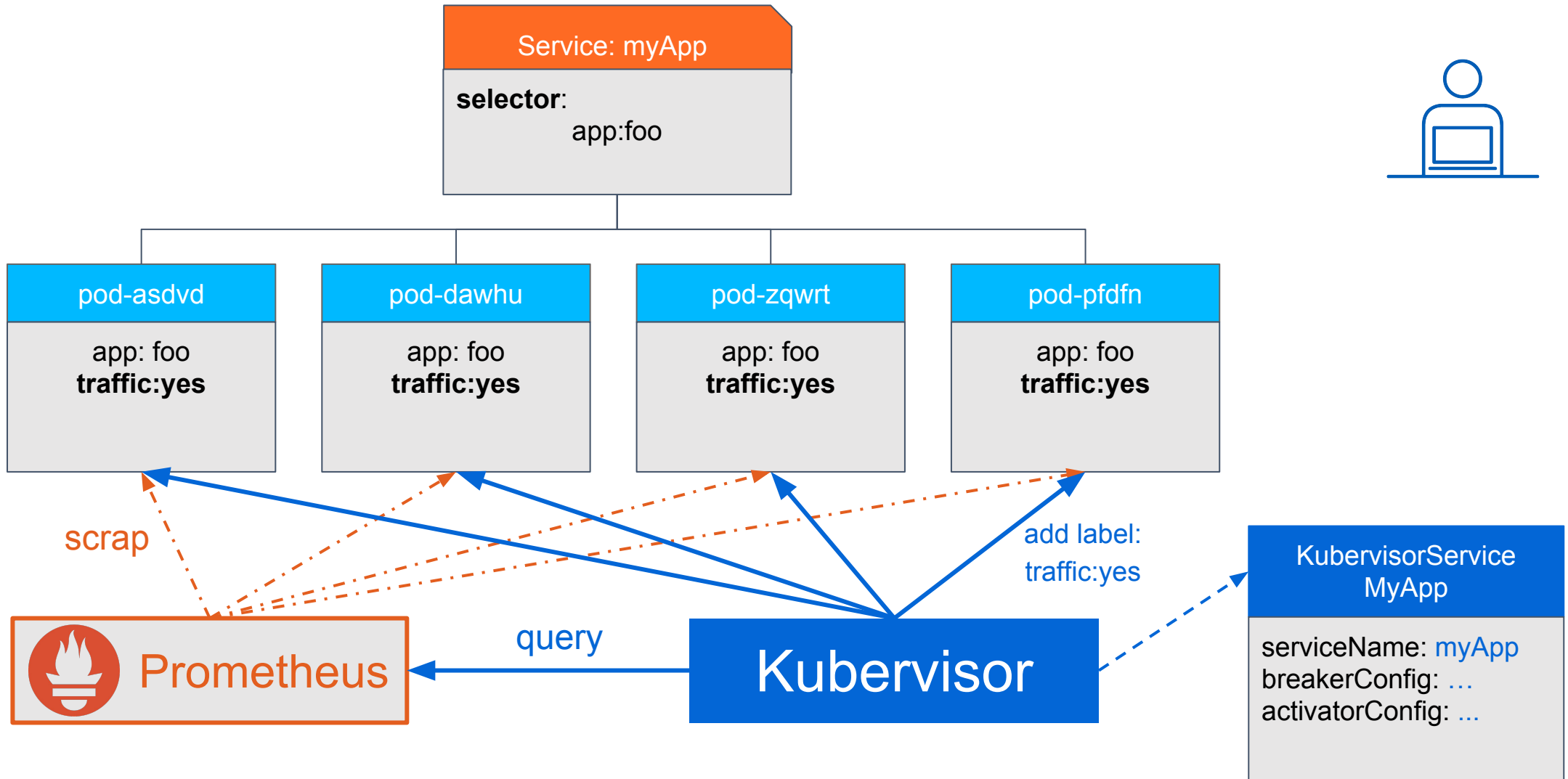


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor workflow

## Initialisation

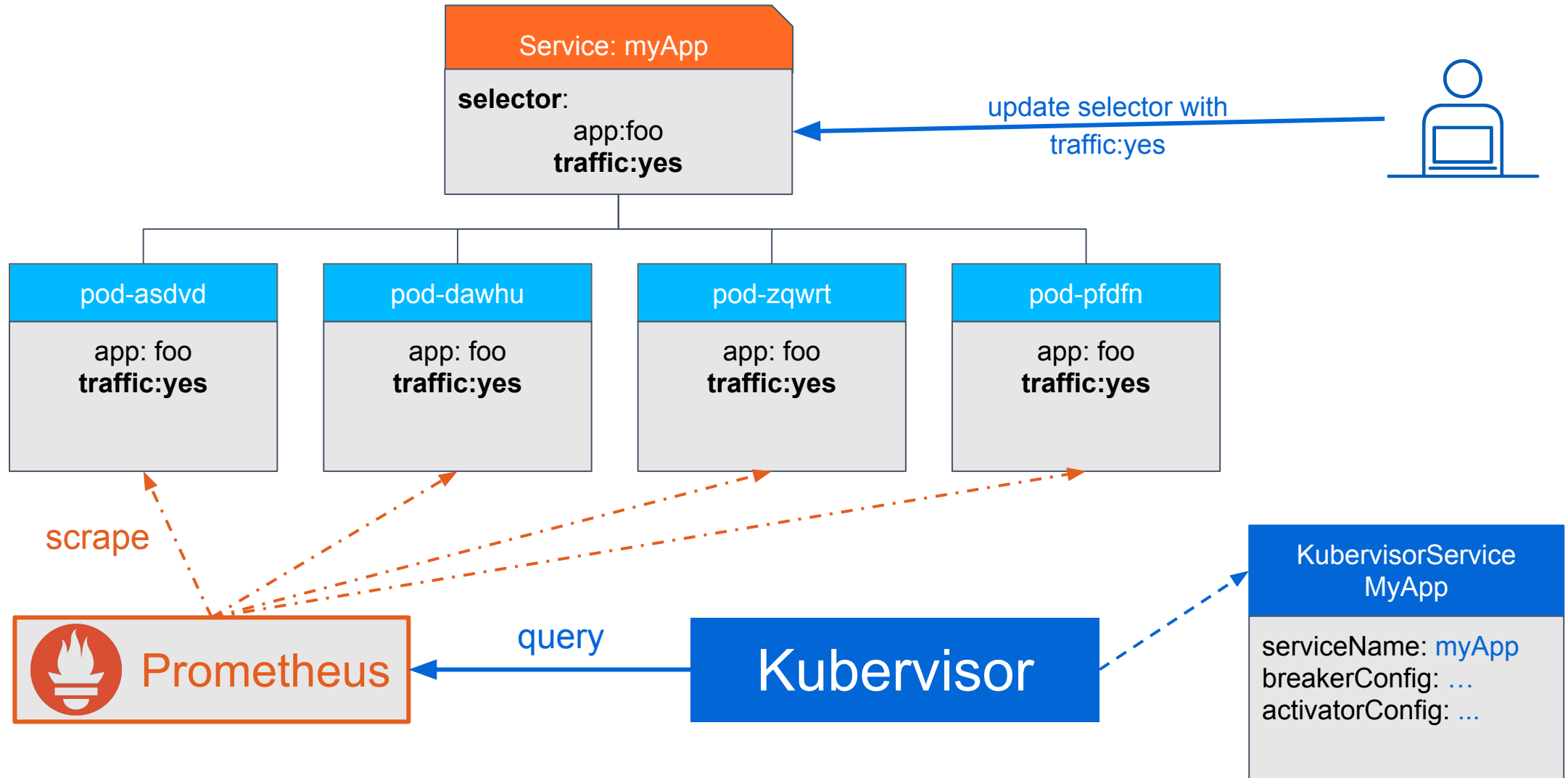


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor workflow

## Anomaly detection

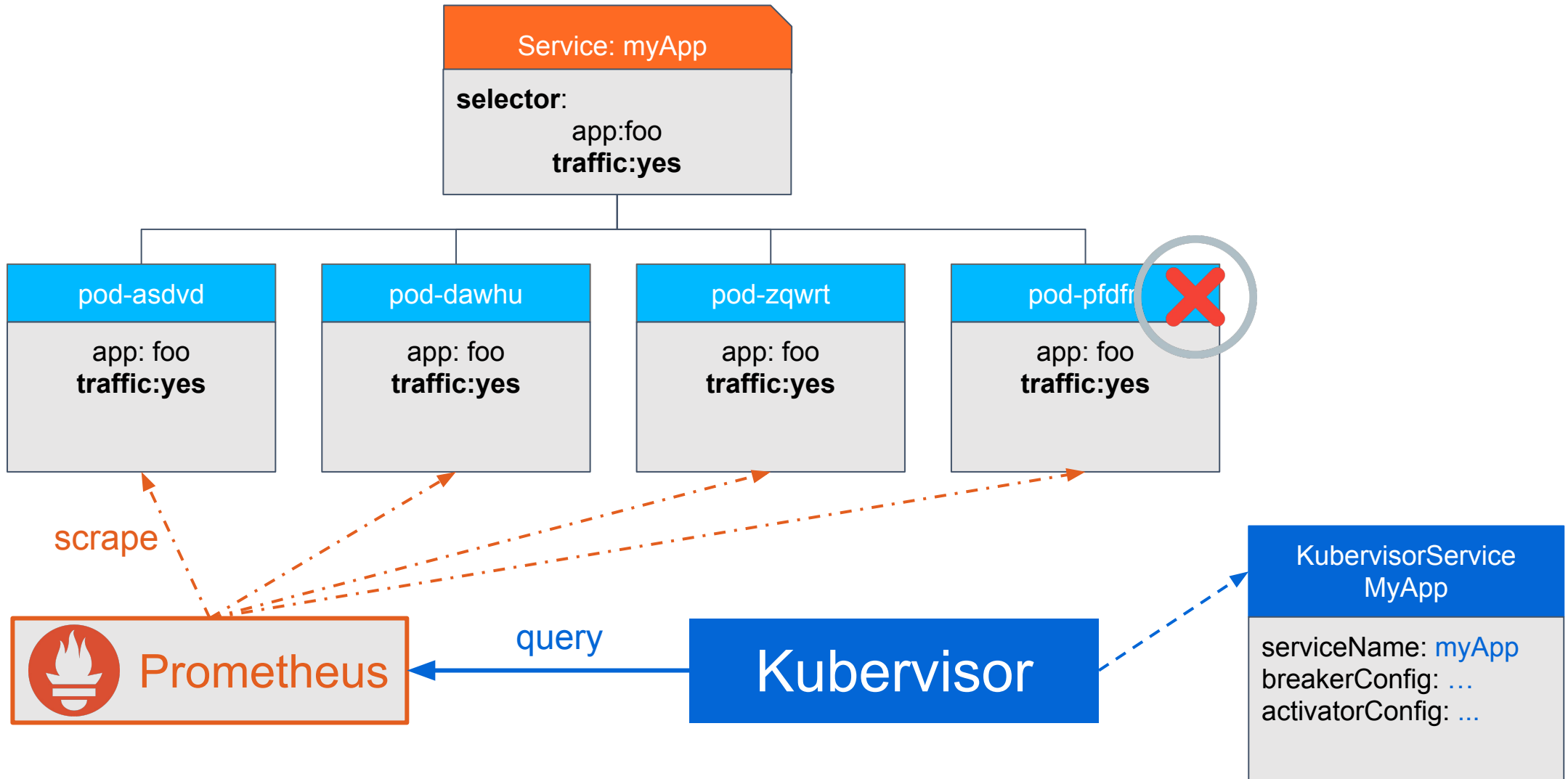


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor workflow

## Breaker activation

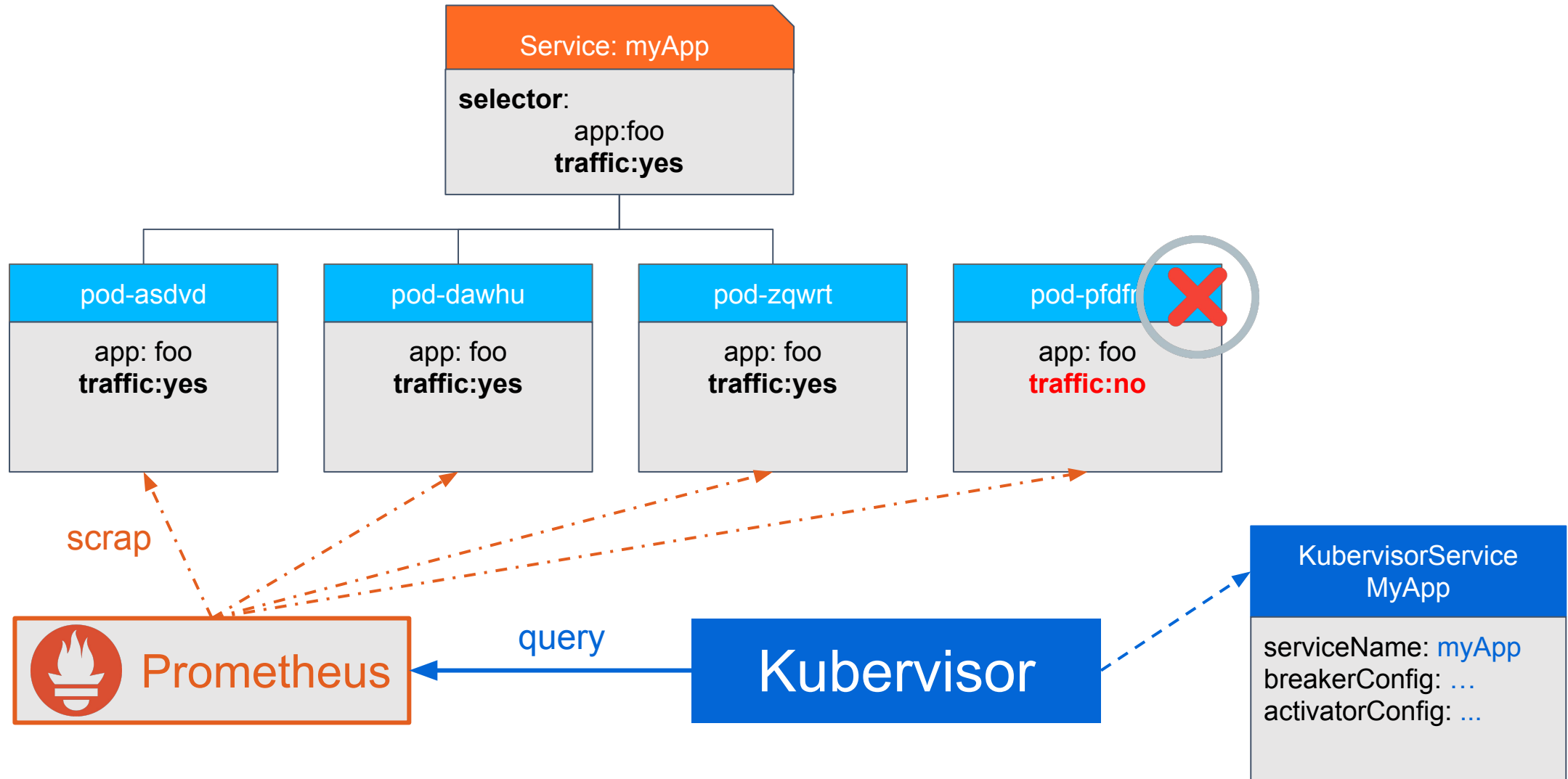


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor

A pod anomaly detection solution



KubeCon



CloudNativeCon

Europe 2018

# Demo time!

# Demo time !



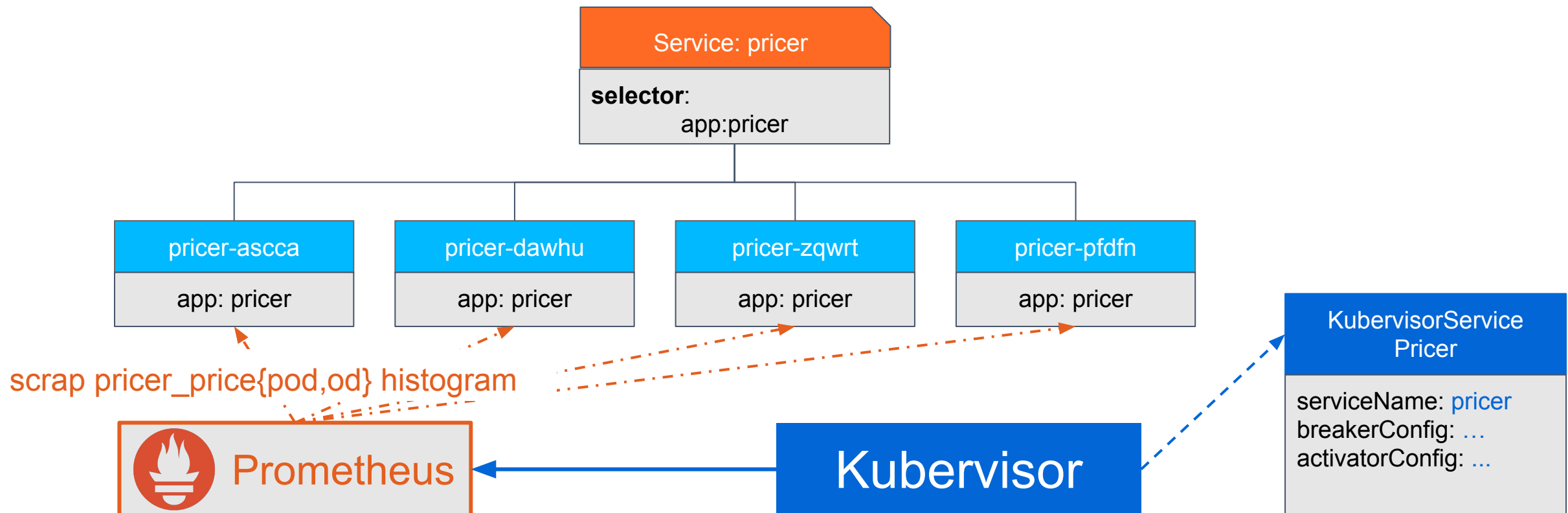
KubeCon



CloudNativeCon

Europe 2018

- **Flight price search:** return 3 best prices for an "origin" and "Destination" city on the current day.





# Kubervisor

## Internal Architecture

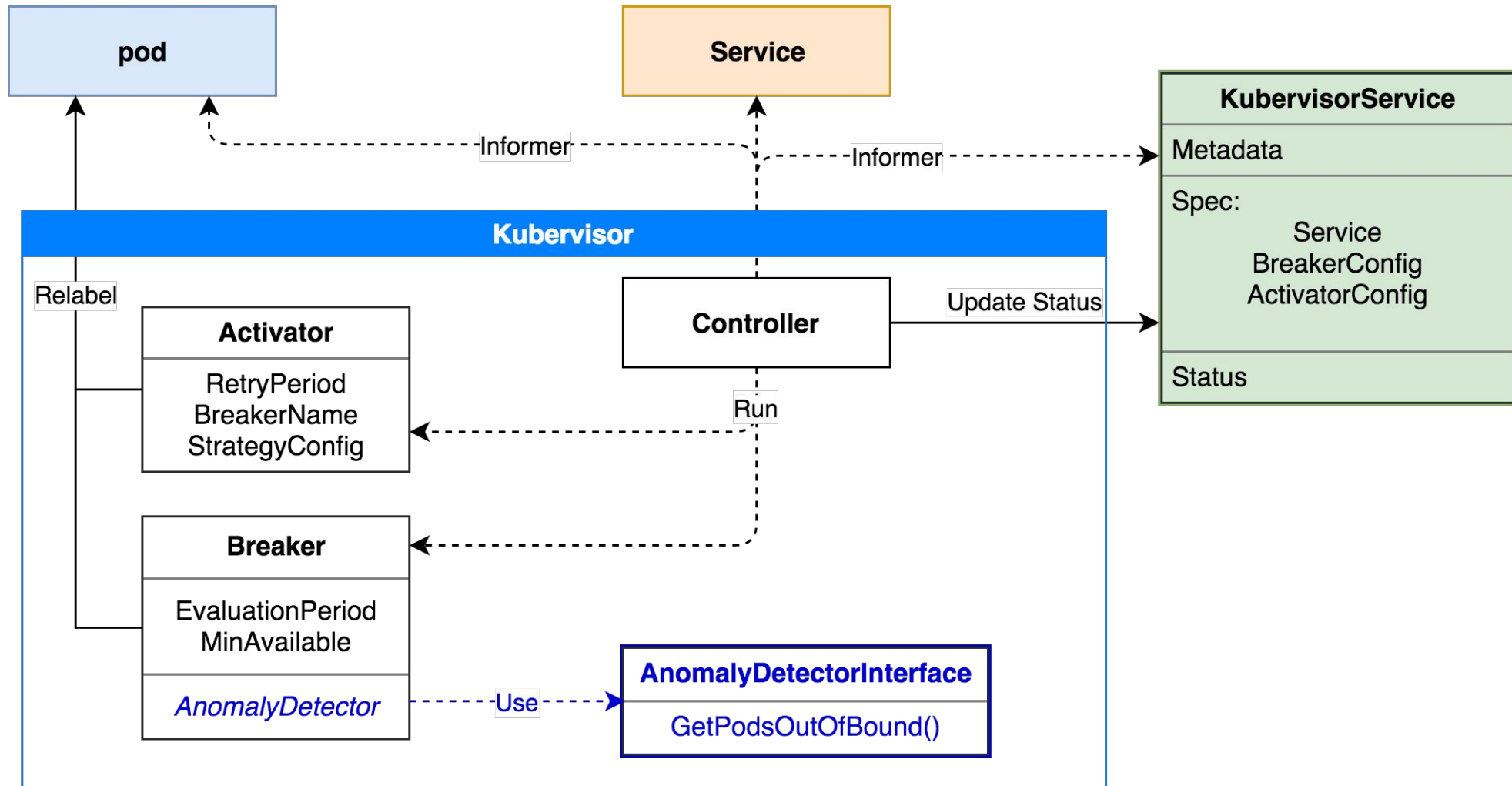


KubeCon



CloudNativeCon

Europe 2018



# Kubervisor

## Internal Architecture

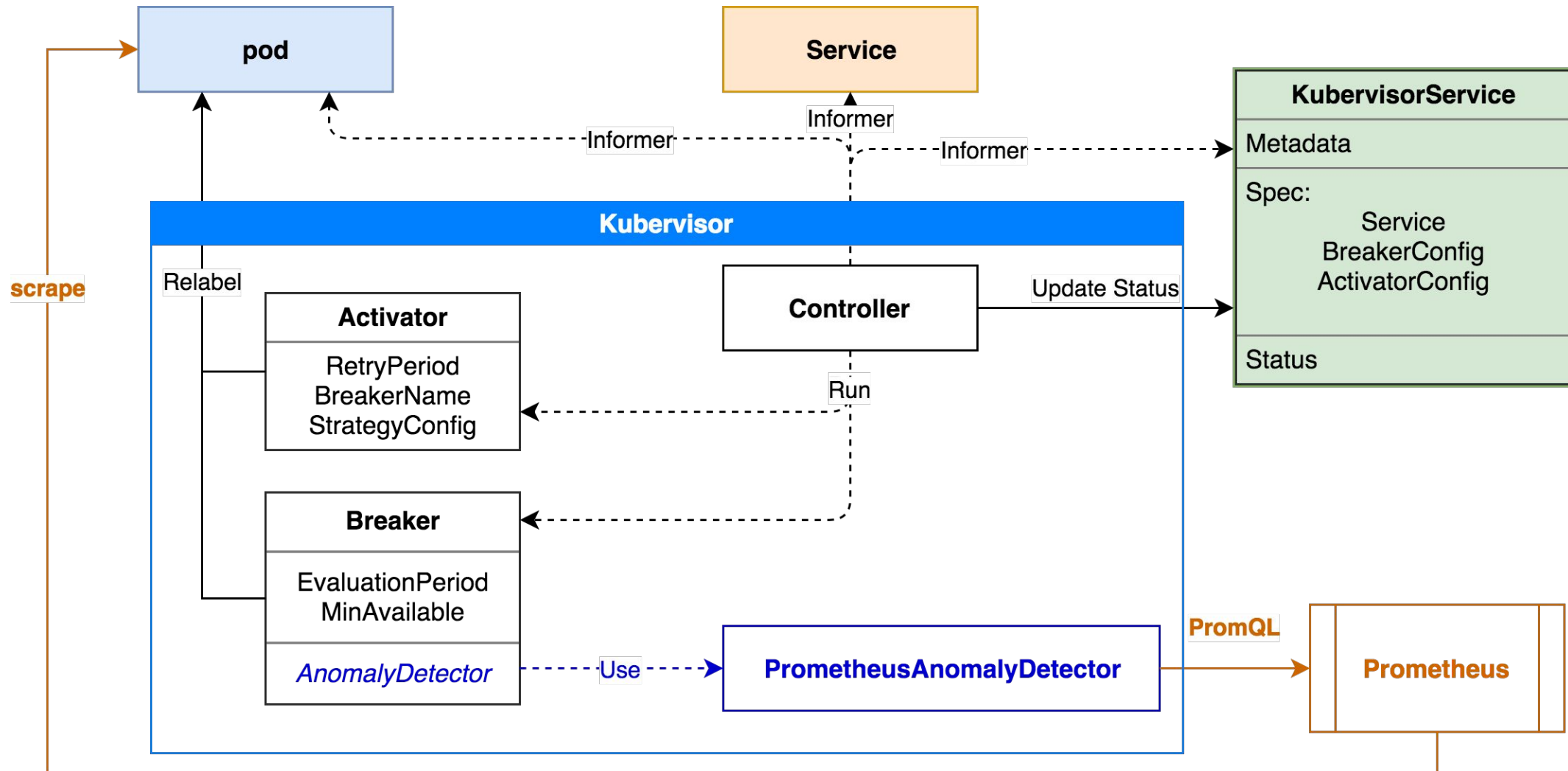


KubeCon



CloudNativeCon

Europe 2018



# Possible Extension



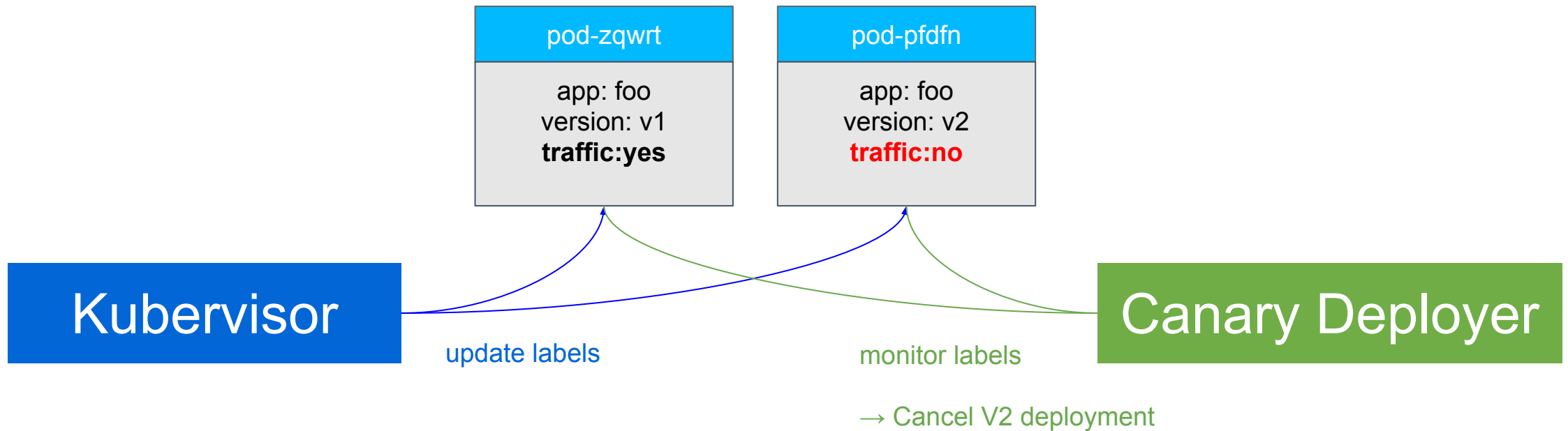
KubeCon



CloudNativeCon

Europe 2018

- Kubervisor labels can be used by other controller
- Example Canary testing deployment



# Key takeaways



KubeCon



CloudNativeCon

Europe 2018

- Additional solution for improving your services reliability
- Based on standard: Controller with CRD, PromQL (Prometheus)
- Extendable: plugable Breaker and Activator implementation.
- Open source: test it, break it and open us issues



KubeCon



CloudNativeCon

Europe 2018

- [github.com/amadeusitgroup/kubervisor](https://github.com/amadeusitgroup/kubervisor)

- @cedriclam

- @BenqueDavid

# Questions?

