



**KubeCon**



**CloudNativeCon**

**North America 2024**





KubeCon



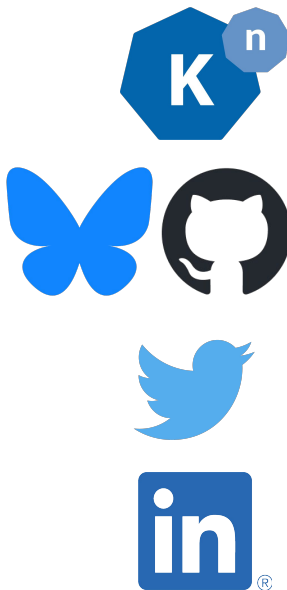
CloudNativeCon

North America 2024

# All your routes are ready, *more or less*

David Protasowski





**Staff Engineer**

**Serving Lead & TOC**

**dprotaso**

**lintinmybelly**

**dprotasowski**

- **Serving**
  - Autoscale your workloads based on traffic - can scale to zero
- **Eventing**
  - Declaratively bind consumers and producers of events
- **Client (kn)**
  - Create resources interactively from the command line
- **Functions**
  - A programming model to simplify development



- **Serving**
  - Autoscale your workloads based on traffic - can scale to zero
- **Eventing**
  - Declaratively bind consumers and producers of events
- **Client (kn)**
  - Create resources interactively from the command line
- **Functions**
  - A programming model to simplify development



- Higher Level Abstraction than Kubernetes
- Container  $\Rightarrow$  URL
- Autoscales your workloads based on traffic
  - Scales to Zero
- Automatic Certificate Provisioning for HTTPS
- Revision Management & Traffic Splitting
- Automatic Health Checks



```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
        - image: my-func:2.0
          ports:
            - containerPort: 80
  traffic:
    - percent: 0
      revisionName: my-func-0004
      tag: staging
    - percent: 40
      revisionName: my-func-0002
    - percent: 60
      revisionName: my-func-0003
```

<https://my-func.default.example.com>

# Knative Serving

```
apiVersion: autoscaling/v2
kind: HorizontalPodAutoscaler
metadata:
  name: my-func
spec:
  selector:
    matchLabels:
      app: my-func
  replicas: 2
  template:
    metadata:
      labels:
        app: my-func
    spec:
      containers:
        - name: my-func
          image: my-func:2.0
          ports:
            - containerPort: 80
```

```
apiVersion: v1
kind: Service
metadata:
  name: my-func
spec:
  selector:
    app: my-func
  ports:
    - port: 80
      protocol: TCP
      targetPort: 80
```

```
apiVersion: v1
kind: Ingress
metadata:
  name: my-func
spec:
  rules:
    - http: my-func
      paths:
        - path: /
          backend:
            serviceName: my-func
            servicePort: 80
```



```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
        - image: my-func:2.0
          ports:
            - containerPort: 80
```

<https://my-func.default.example.com>



# Knative Serving



KubeCon



CloudNativeCon

North America 2024

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
```

# Knative Serving



KubeCon



CloudNativeCon

North America 2024

**Service**

# Knative Serving



KubeCon



CloudNativeCon

North America 2024

**Service**

# Knative Serving

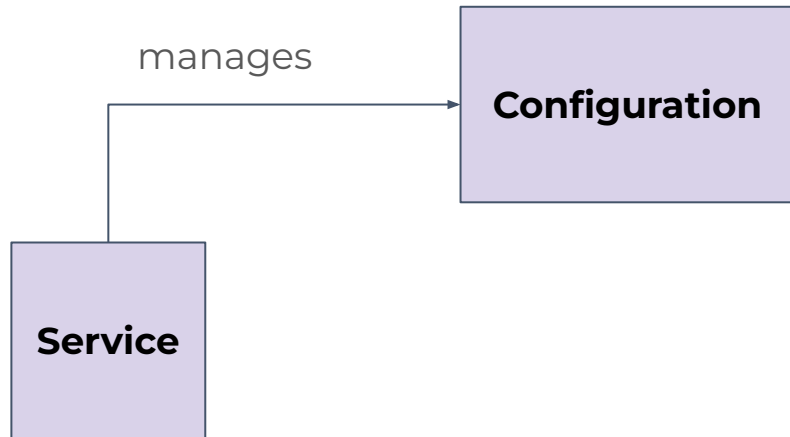


KubeCon

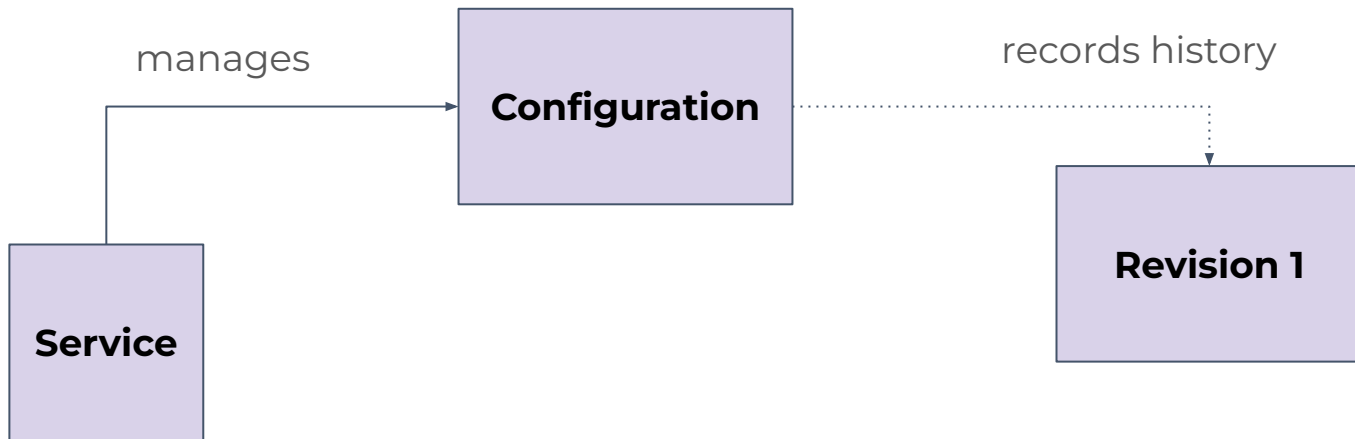


CloudNativeCon

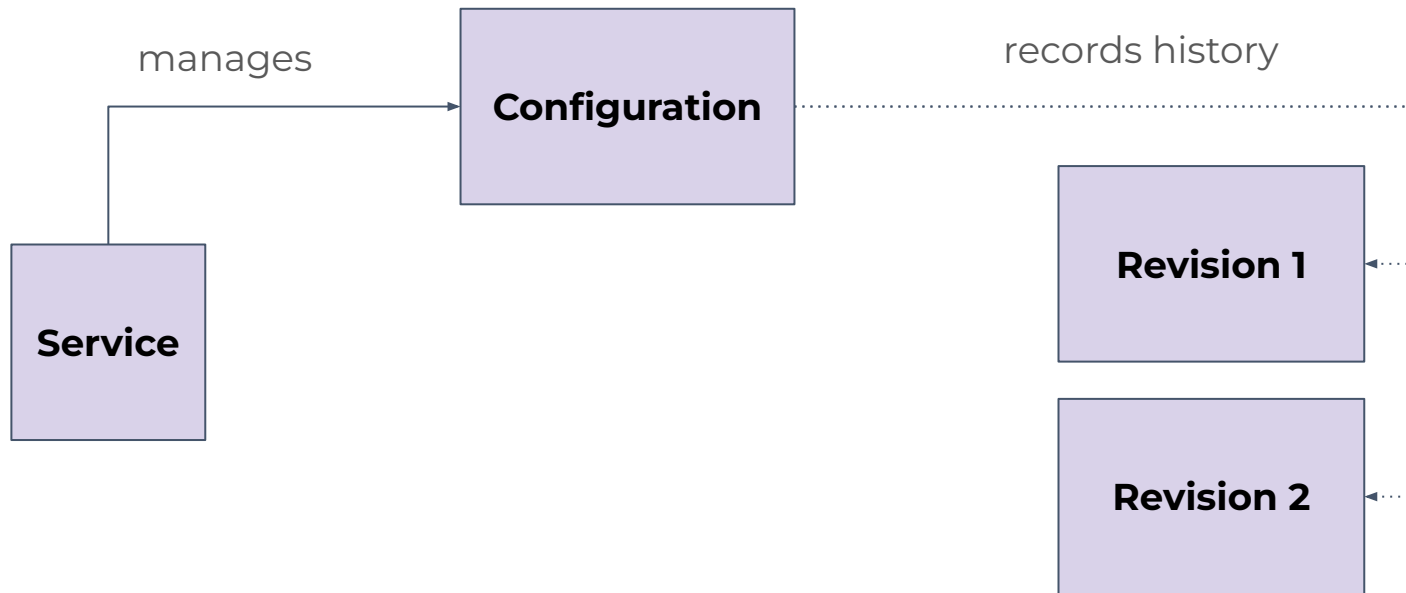
North America 2024



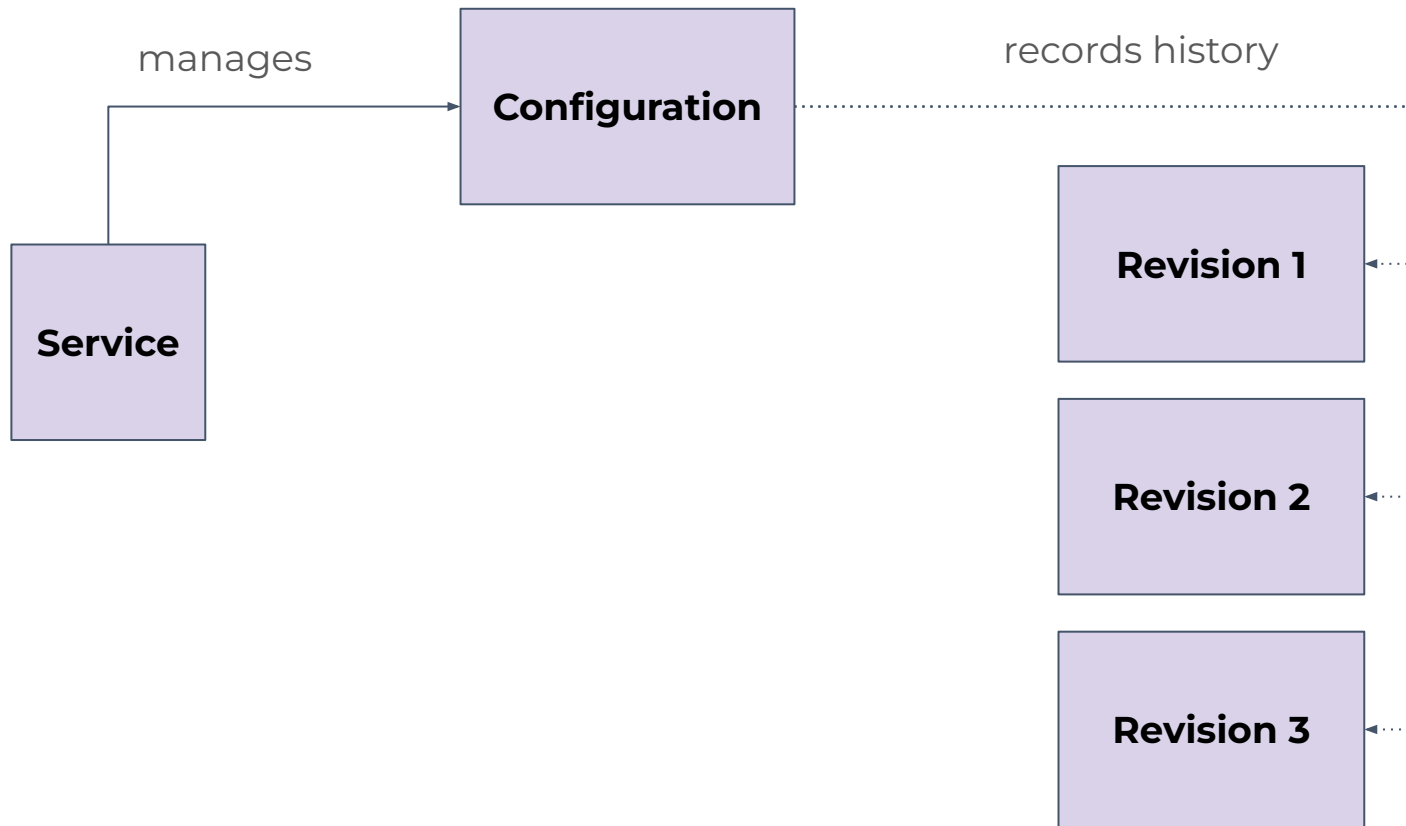
# Knative Serving



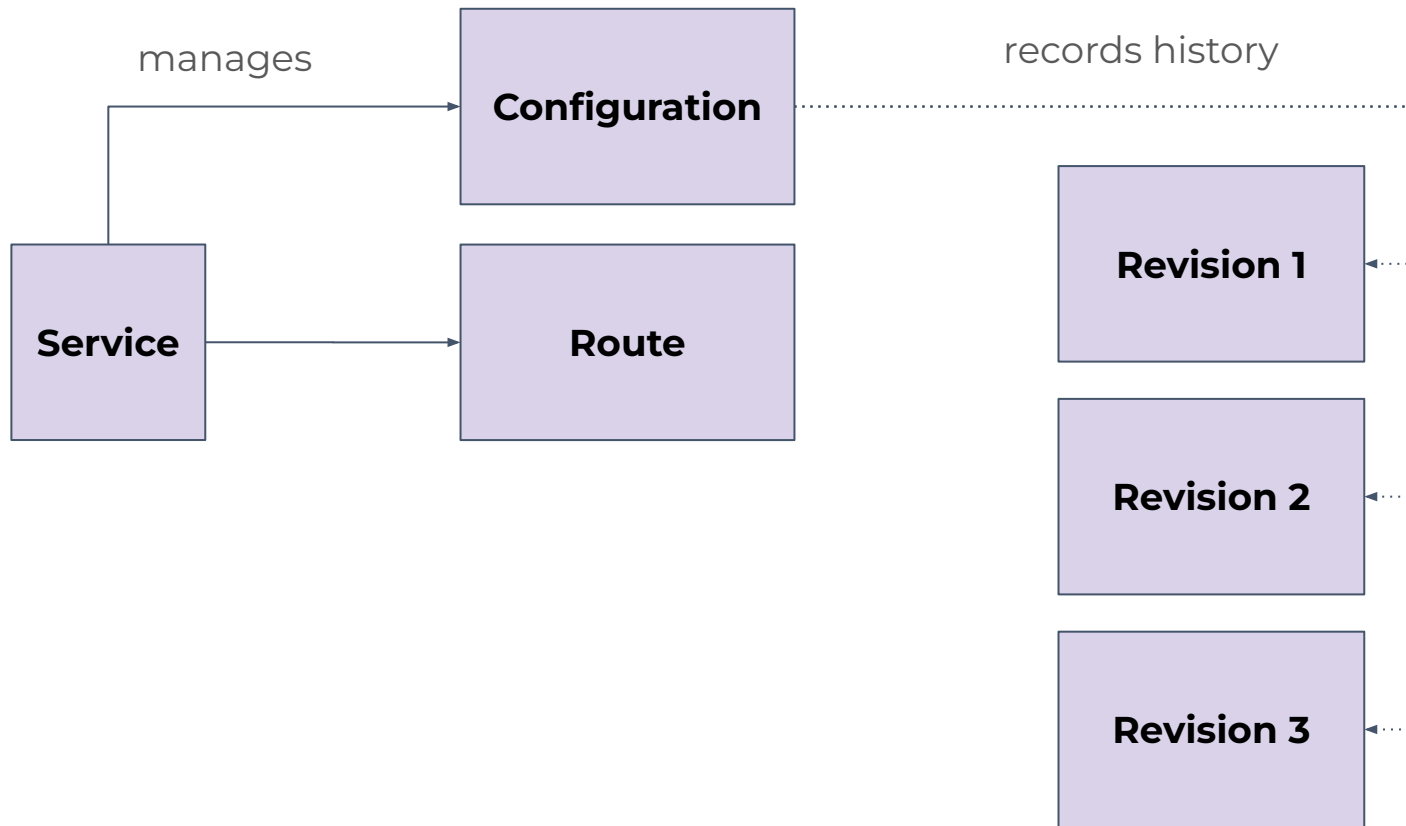
# Knative Serving



# Knative Serving

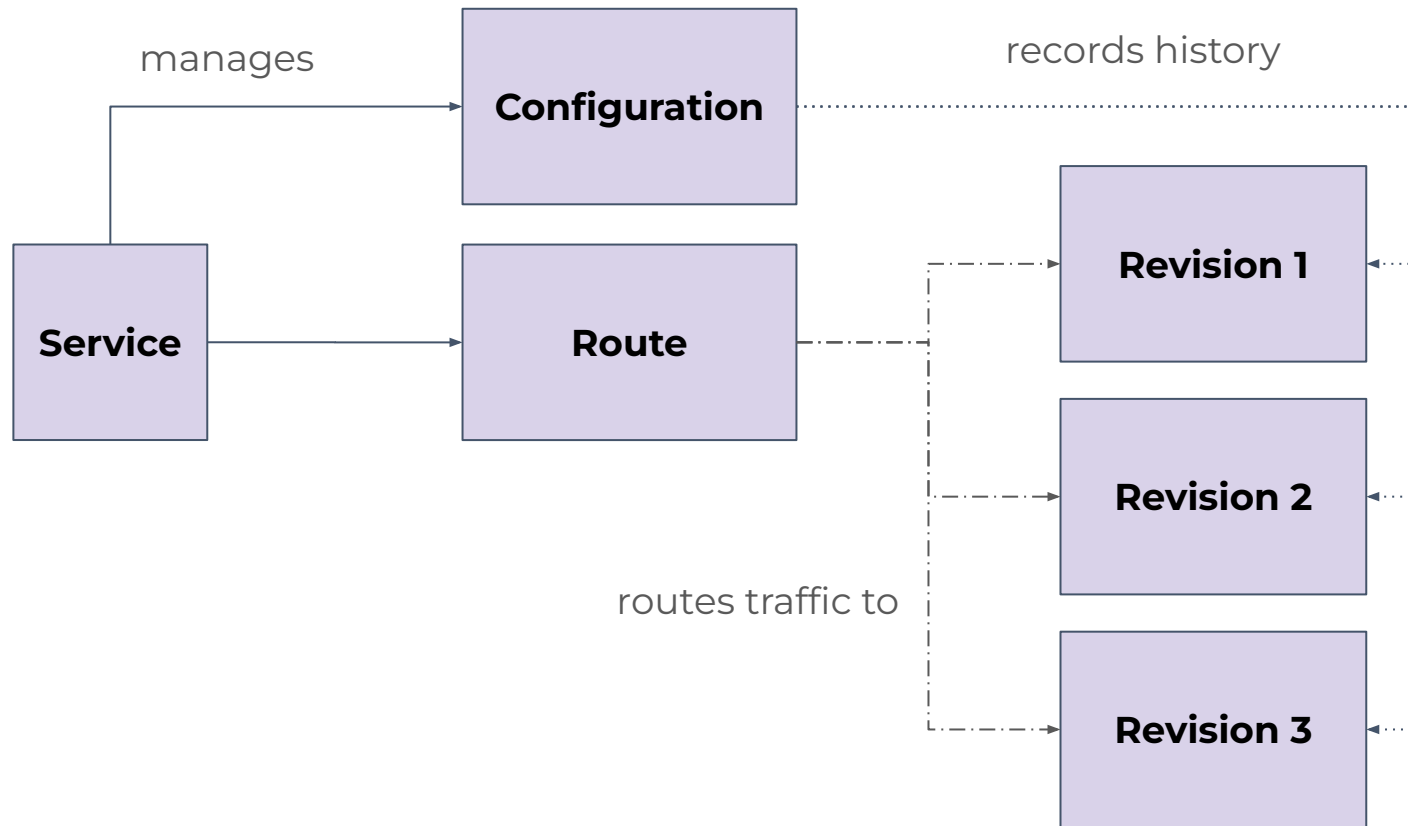


# Knative Serving





# Knative Serving



# Knative Serving - Routes

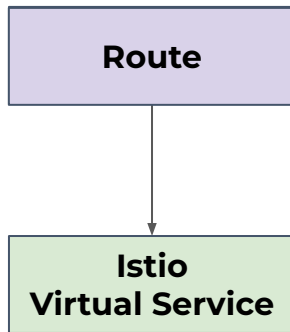


KubeCon



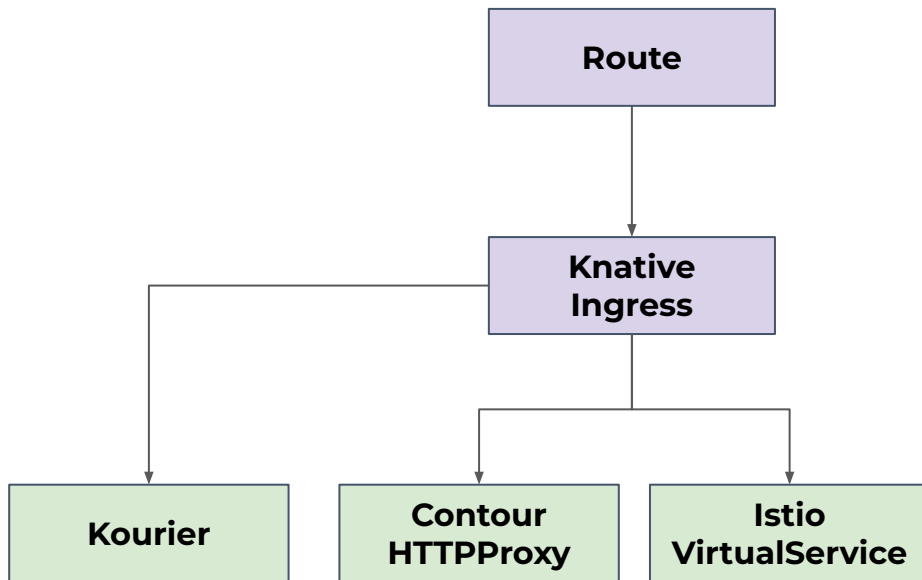
CloudNativeCon

North America 2024



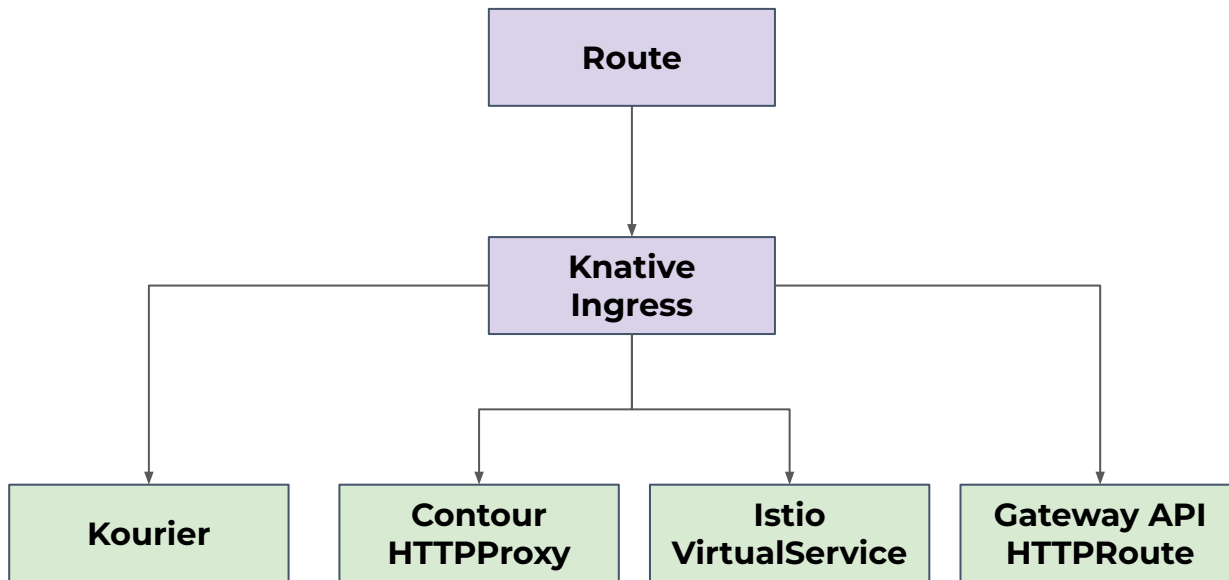
Q1 2018

# Knative Serving - Routes



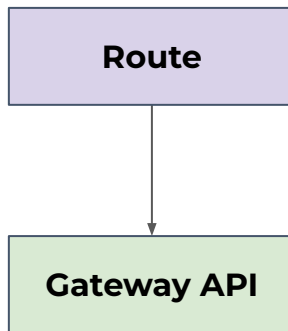
Q4 2018

# Knative Serving - Routes



Q4 2021

# Knative Serving - Routes



**Future**



KubeCon



CloudNativeCon

North America 2024

# What have we learned?



KubeCon



CloudNativeCon

North America 2024

All your routes are ready,  
*more or less*



KubeCon



CloudNativeCon

North America 2024

All your routes are **ready**,  
*more or less*



# Knative Serving - Ready Condition

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
```

# Knative Serving - Ready Condition



KubeCon



CloudNativeCon

North America 2024

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
status:
  ...
  conditions: []
```

# Knative Serving - Ready Condition



KubeCon



CloudNativeCon

North America 2024

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
status:
  ...
  conditions:
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: ConfigurationsReady
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: RoutesReady
```

# Knative Serving - Ready Condition

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
status:
  ...
  conditions:
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: ConfigurationsReady
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: RoutesReady
```



# Knative Serving - Ready Condition



KubeCon



CloudNativeCon

North America 2024

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
status:
  ...
  conditions:
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: ConfigurationsReady

  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: RoutesReady

  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: Ready
```

# Knative Serving - Ready Condition

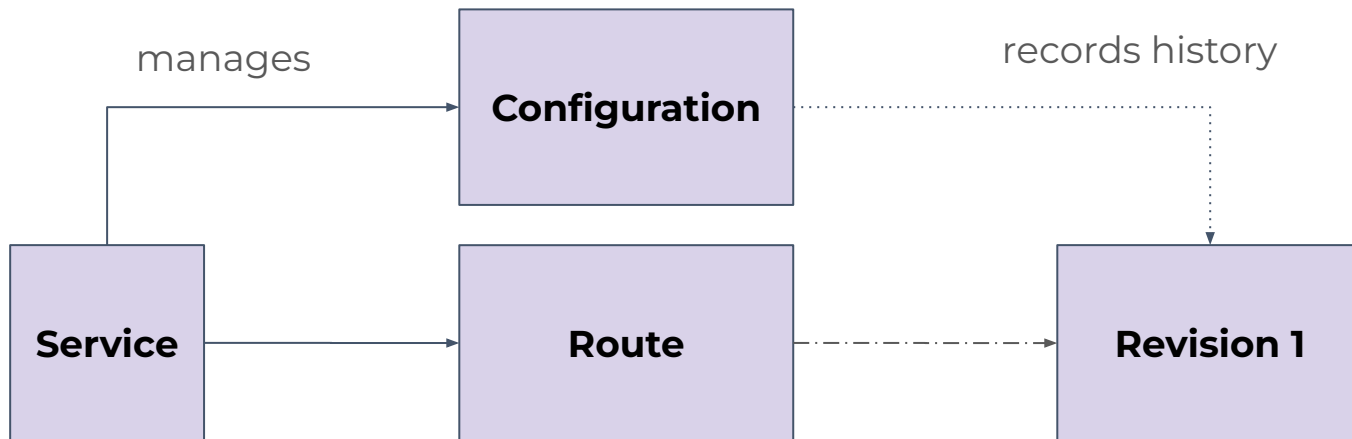
```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: my-func
spec:
  template:
    spec:
      containers:
      - image: my-func:2.0
        ports:
        - containerPort: 80
status:
  ...
  conditions:
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: ConfigurationsReady

  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: RoutesReady

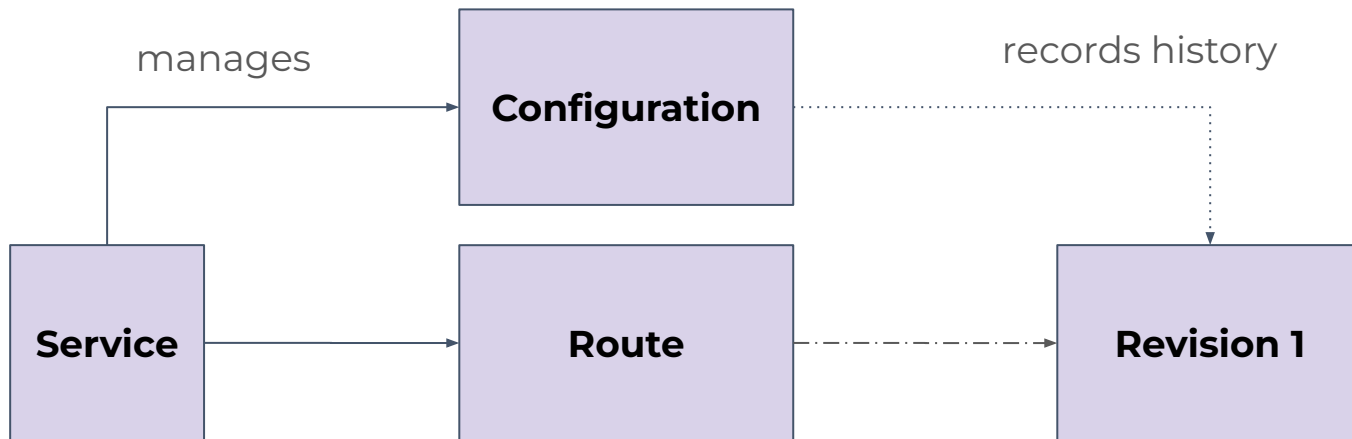
  - lastTransitionTime: "2024-11-03T18:20:01Z"
    status: "True"
    type: Ready
```



# Knative Serving - Ready Condition

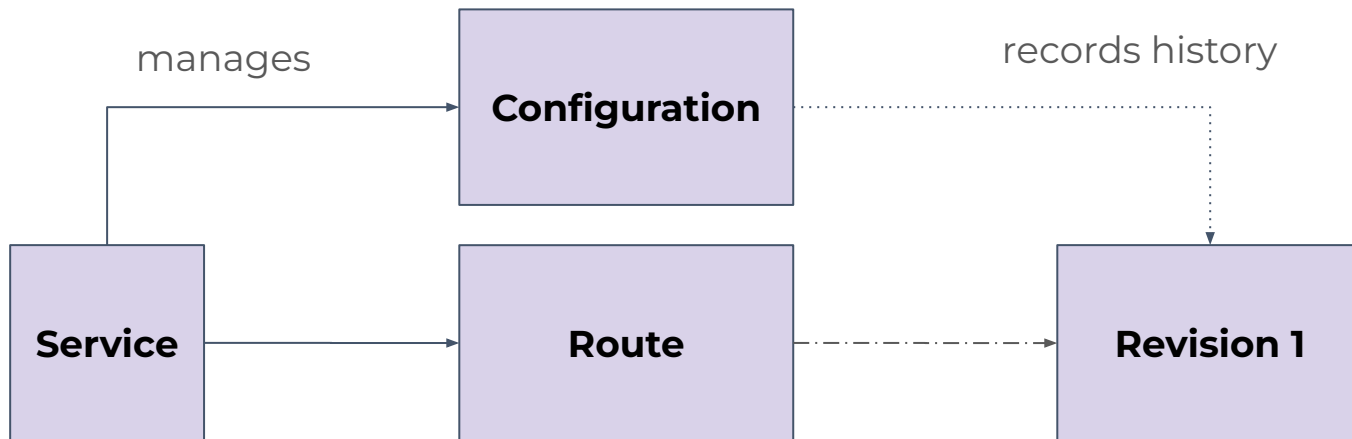


# Knative Serving - Ready Condition

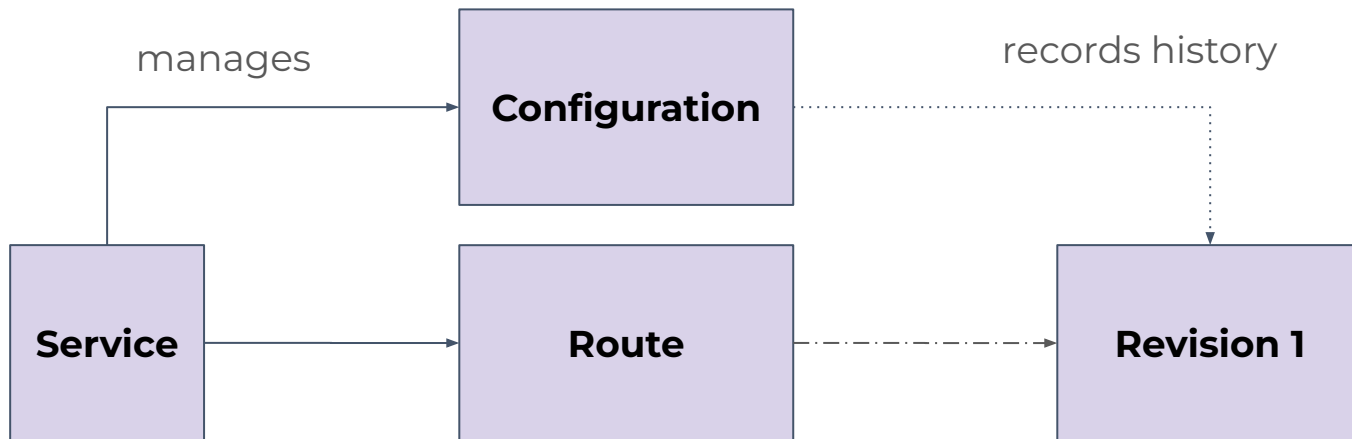




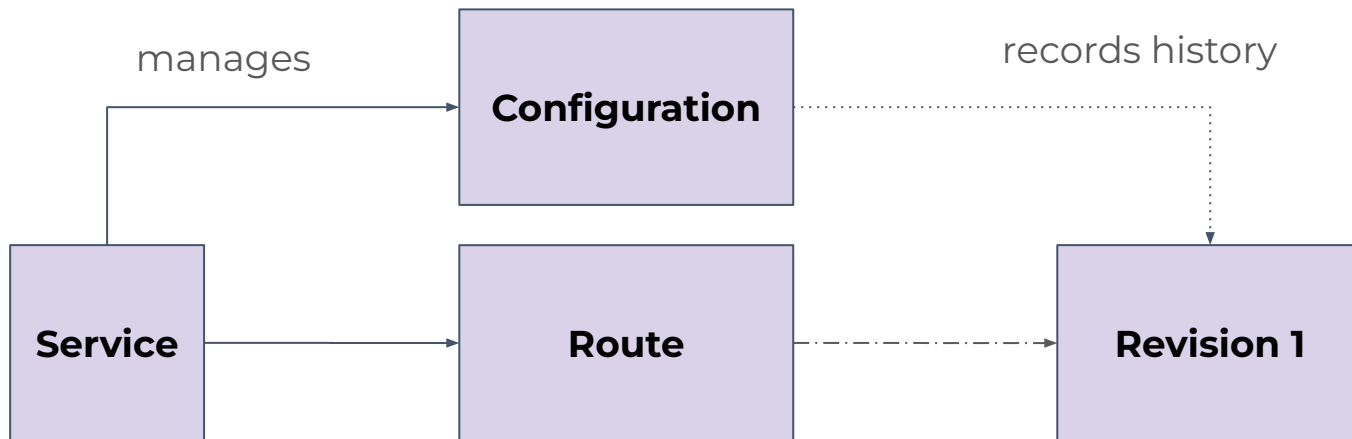
# Knative Serving - Ready Condition



# Knative Serving - Ready Condition



# Knative Serving - Ready Condition





KubeCon



CloudNativeCon

North America 2024

All your routes are **ready**,  
*more or less*



KubeCon



CloudNativeCon

North America 2024

All your **routes** are ready,  
*more or less*

# Knative Serving - Routes

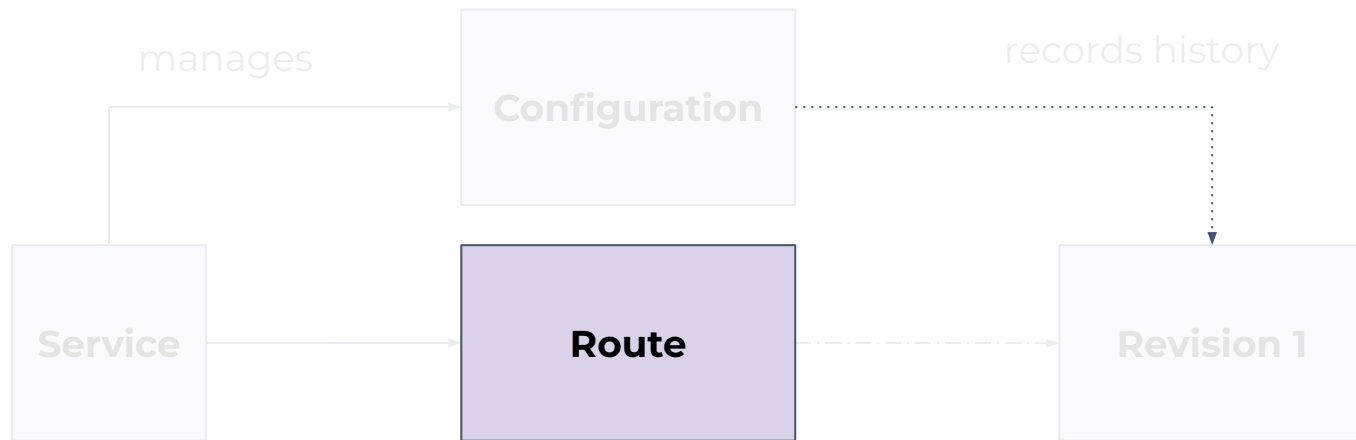


KubeCon

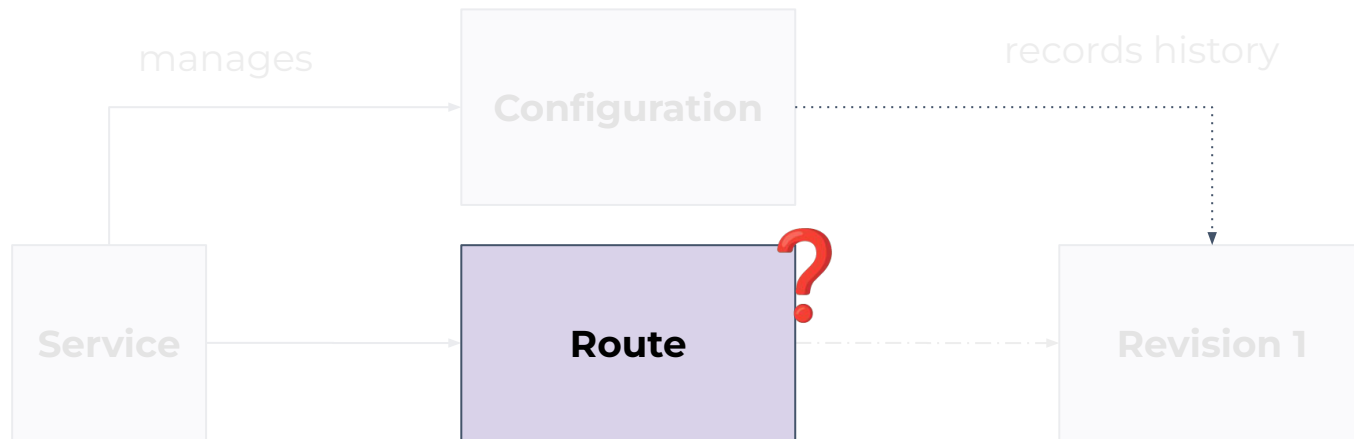


CloudNativeCon

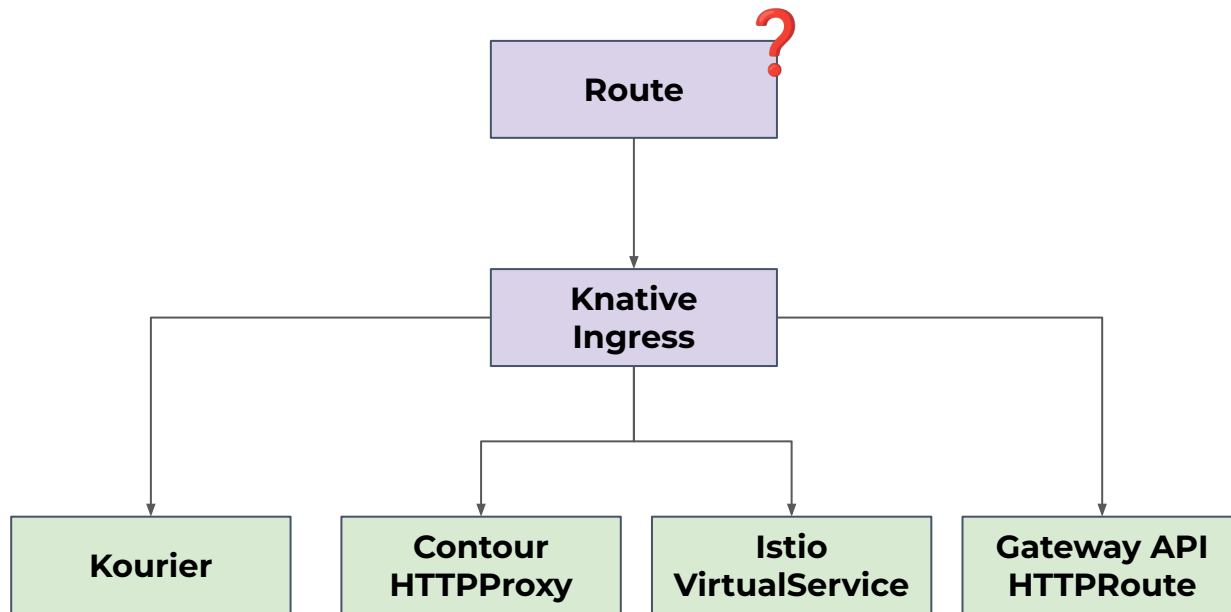
North America 2024



# Knative Serving - Routes

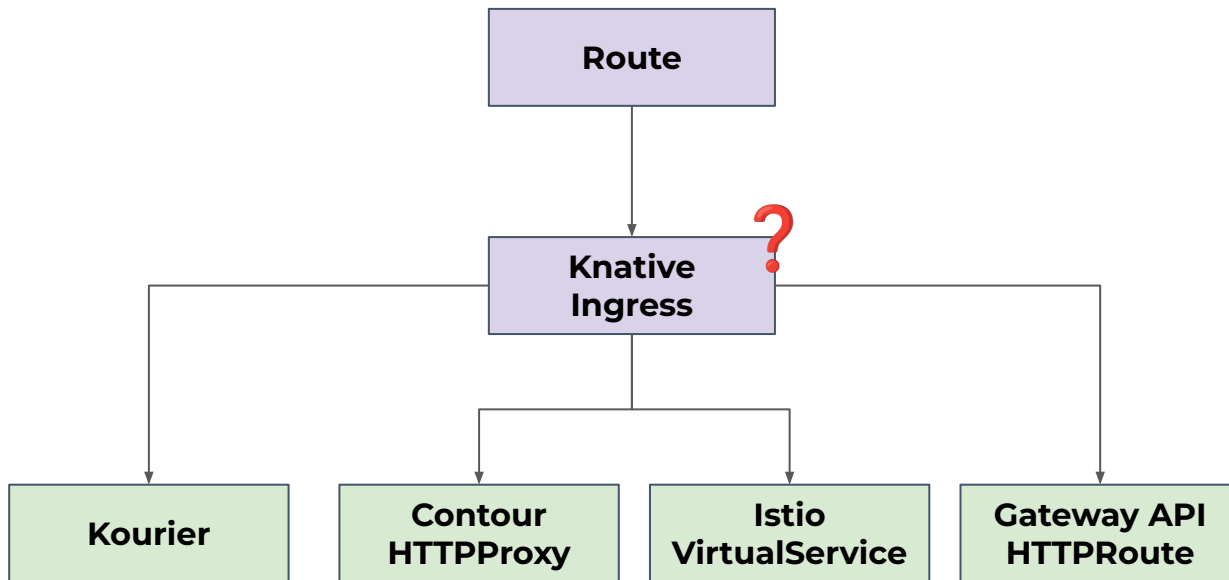


# Knative Serving - Routes

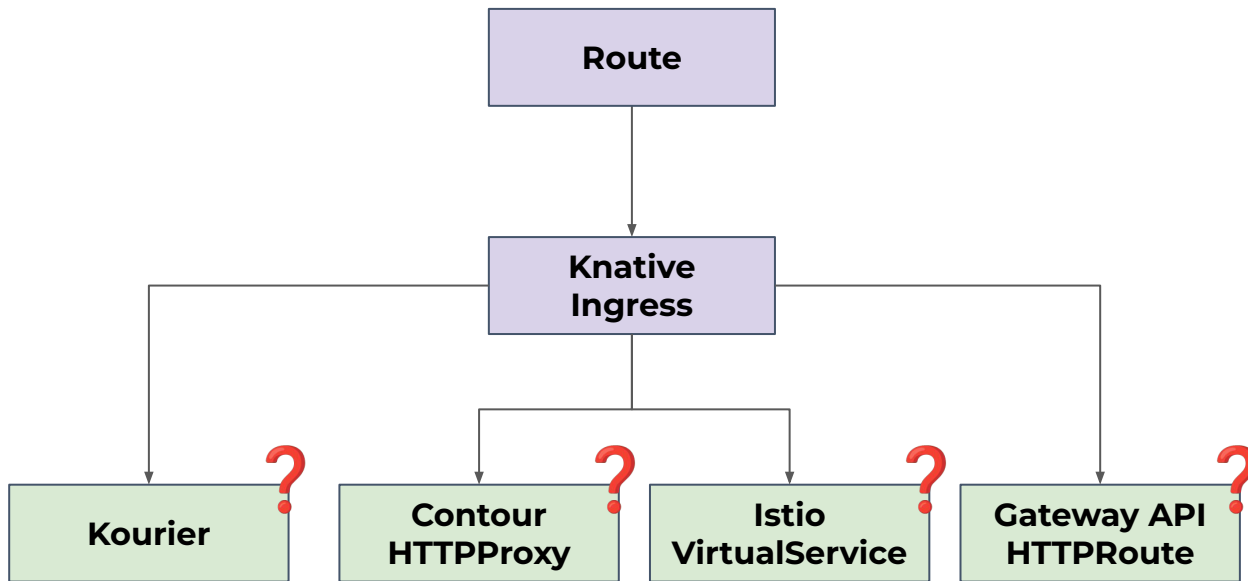




# Knative Serving - Routes



# Knative Serving - Routes



# Types of Guarantees

- Controller validates configuration
- Controllers sends configuration to proxy

# Is it safe?



KubeCon



CloudNativeCon

North America 2024



# Types of Guarantees

- ✗ Controller validates configuration
- ✗ Controllers sends configuration to proxy

- ✗ Controller validates configuration
- ✗ Controllers sends configuration to proxy
  - Proxy responds to config message with an ACK

- ✗ Controller validates configuration
- ✗ Controllers sends configuration to proxy
- ✗ Proxy responds to config message with an ACK

- ✗ Controller validates configuration
- ✗ Controllers sends configuration to proxy
- ✗ Proxy responds to config message with an ACK
  - Proxy waits for programming to finish then replies (maybe async)



- ✗ Controller validates configuration
- ✗ Controllers sends configuration to proxy
- ✗ Proxy responds to config message with an ACK
- ✗ Proxy waits for programming to finish then replies (maybe async)

# DNS

- ✗ Controller validates configuration
- ✗ Controller sends configuration to proxy
- ✗ Proxy responds to config message with ACK
- ✗ Proxy waits for programming to finish then replies (maybe async)

>1



Controller validates configuration



Controller sends configuration to proxy



Proxy responds to config message with an ACK



Proxy starts programming to fulfill the replies  
(maybe asynchronously)

# PROXY

- **Accepted**

- If the Route's ParentRef specifies an existing Gateway that supports Routes of this kind AND that Gateway's controller has sufficient access, then that Gateway's controller MUST set the "Accepted" condition on the Route

- **Programmed**

- This condition indicates whether a Listener has generated some configuration that will soon be ready in the underlying data plane.

# NOT SAFE



KubeCon



CloudNativeCon

North America 2024



- Some non-Gateway API implementations surface readiness via **status**
- Sometimes it's broken
  - ie. **observedGeneration** is missing or not bumped properly
- tl;dr no reliable way for clients to know when to make a request

- **Ready**

- If used in the future, “Ready” will represent the final state where all configuration is confirmed good *and has completely propagated to the data plane.*

## Related Issues

<https://github.com/kubernetes-sigs/gateway-api/issues/1156>  
<https://github.com/kubernetes-sigs/gateway-api/issues/1364>  
<https://github.com/kubernetes-sigs/gateway-api/issues/1870>

# Come help!



KubeCon



CloudNativeCon

North America 2024

## Vote for 'Ready' in Gateway v1.3!

<https://bit.ly/gateway-ready>







KubeCon



CloudNativeCon

North America 2024



# Probing 🤔



KubeCon



CloudNativeCon

North America 2024





**KIngress**

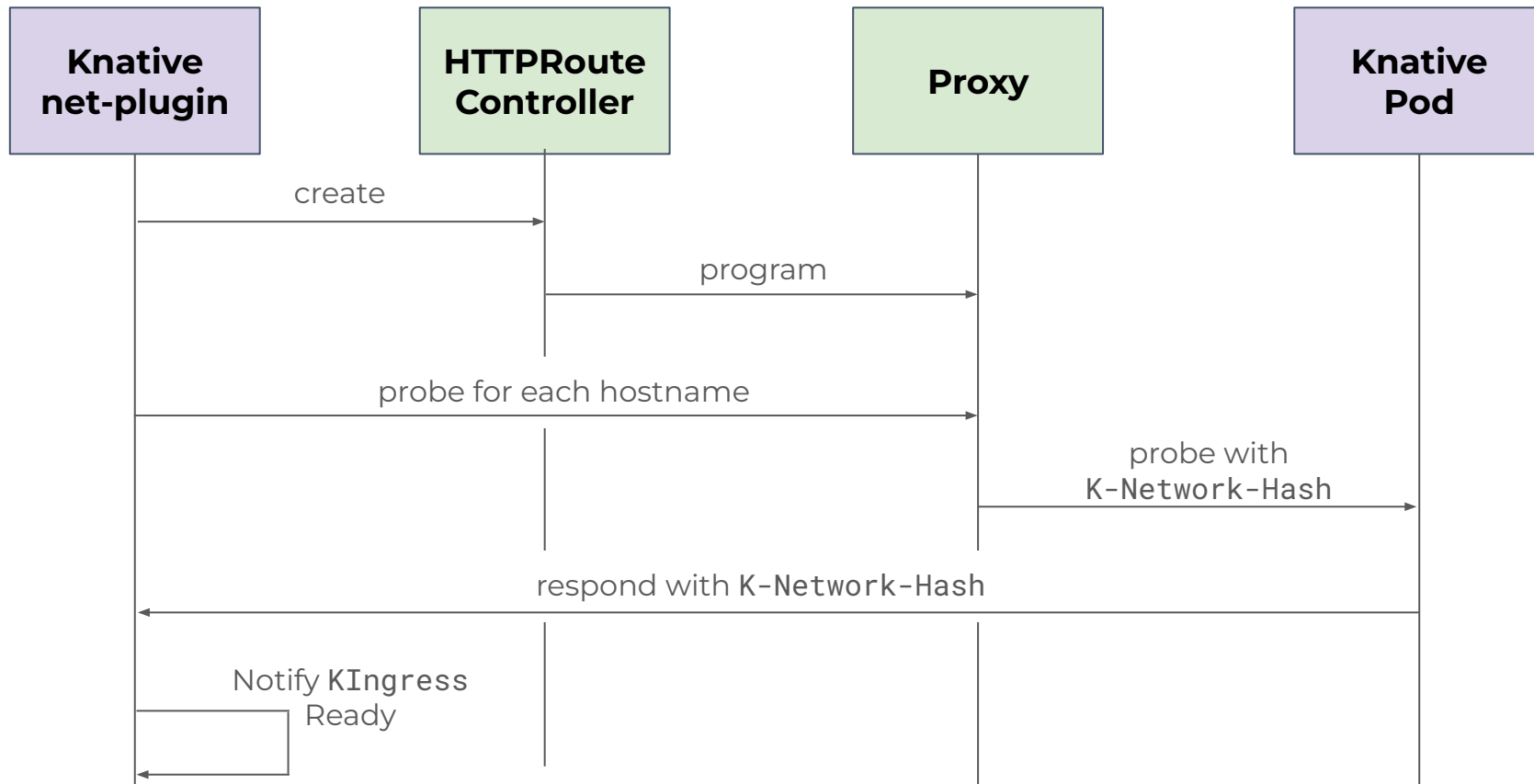
```
apiVersion: gateway.networking.k8s.io/v1
kind: HTTPRoute
metadata:
  name: my-func-route
spec:
  parentRefs:
  - name: knative-gateway
  hostnames:
  - my-func.default.example.com
  rules:
  - backendRefs:
    - name: my-func-0001
      port: 8080
```

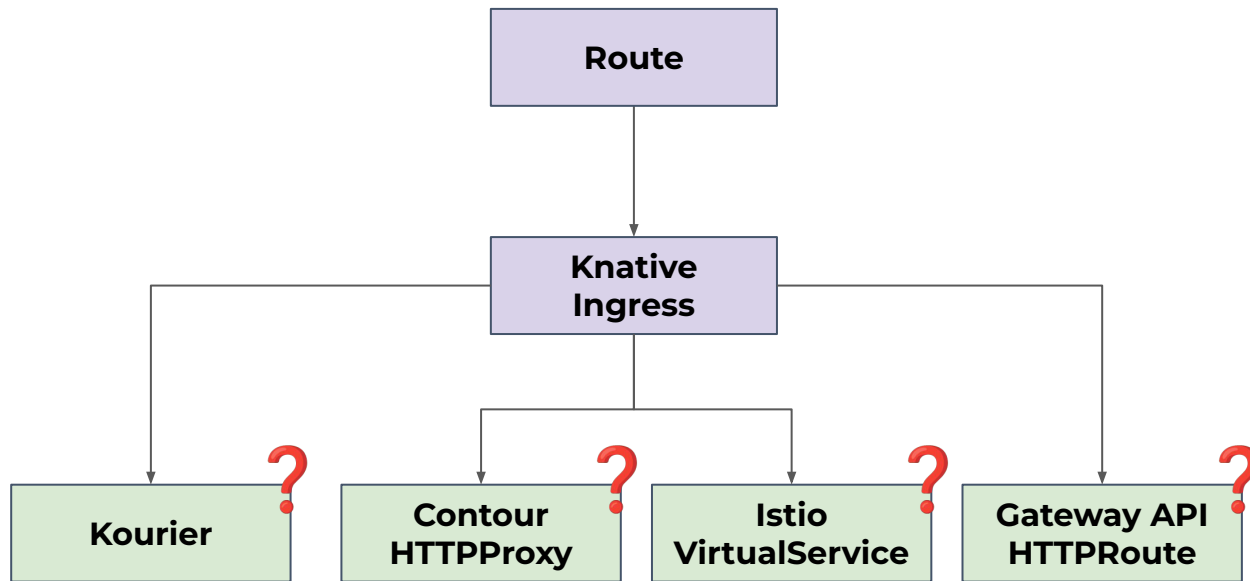


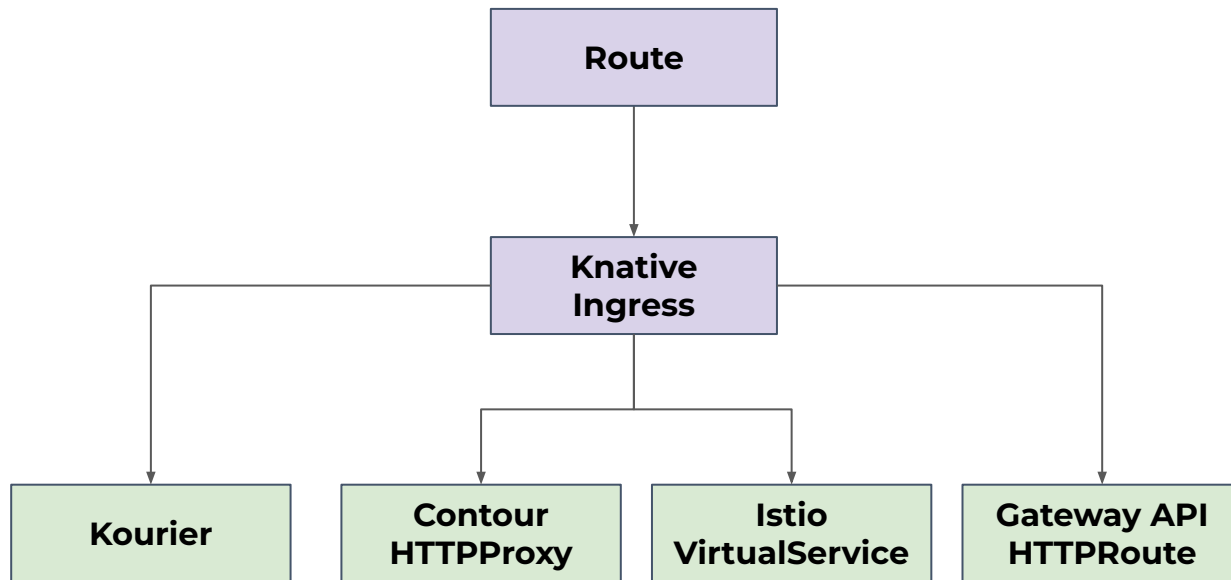
**KIngress**

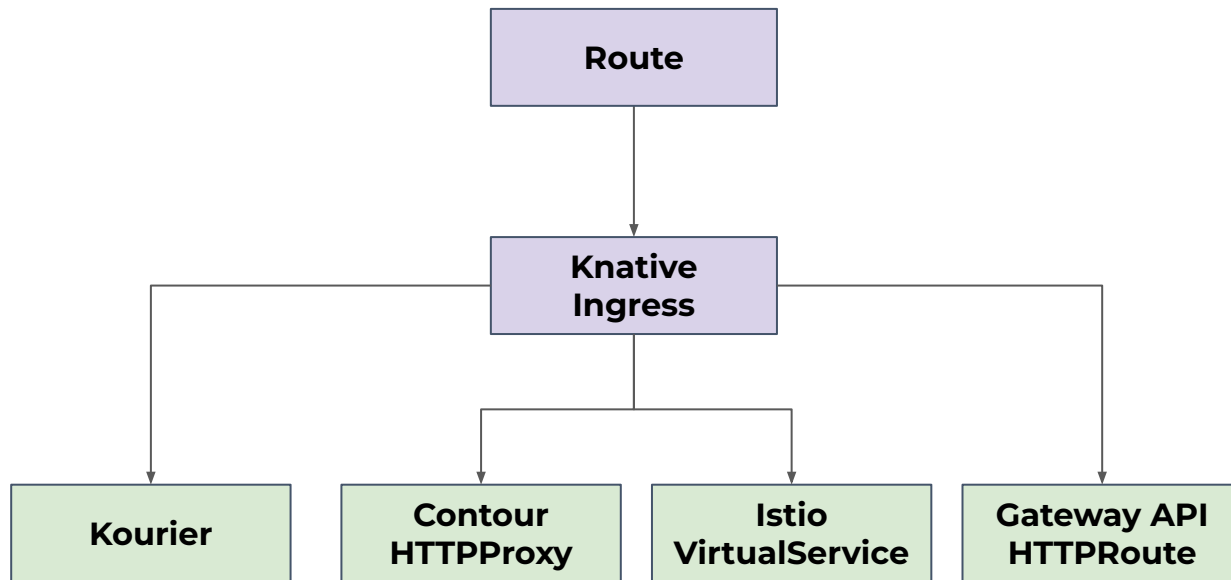
```
apiVersion: gateway.networking.k8s.io/v1
kind: HTTPRoute
metadata:
  name: my-func-route
spec:
  parentRefs:
  - name: knative-gateway
  hostnames:
  - my-func.default.example.com
  rules:
  - matches:
    - headers:
      - type: Exact
        name: K-Network-Hash
        value: override
    filters:
    - type: RequestHeaderModifier
      requestHeaderModifier:
        add:
        - name: K-Network-Hash
          value: a6cf8611a601567cf3b94aba628...
  backendRefs:
  - name: my-func-0001
    port: 8080
```

# Probing 🥲











# We ain't done

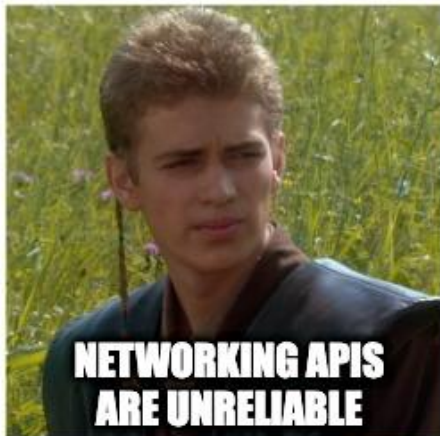


KubeCon



CloudNativeCon

North America 2024



# Update Operations

```
apiVersion: gateway.networking.k8s.io/v1
kind: HTTPRoute
metadata:
  name: my-func-route
spec:
  parentRefs:
  - name: knative-gateway
  hostnames:
  - my-func.default.example.com
  rules:
  - backendRefs:
    - name: my-func-0001
      port: 8080
```



```
apiVersion: gateway.networking.k8s.io/v1
kind: HTTPRoute
metadata:
  name: my-func-route
spec:
  parentRefs:
  - name: knative-gateway
  hostnames:
  - my-func.default.example.com
  rules:
  - backendRefs:
    - name: my-func-0002
      port: 8080
```

```
go run github.com/rakyll/hey@latest \
  -disable-keepalive \
  -c 10 `# number of workers` \
  -q 5 `# qps per worker` \
  -z 1m `# duration` \
  https://my-func.default.example.com
```

# Update Operations are Unsafe



KubeCon



CloudNativeCon

North America 2024

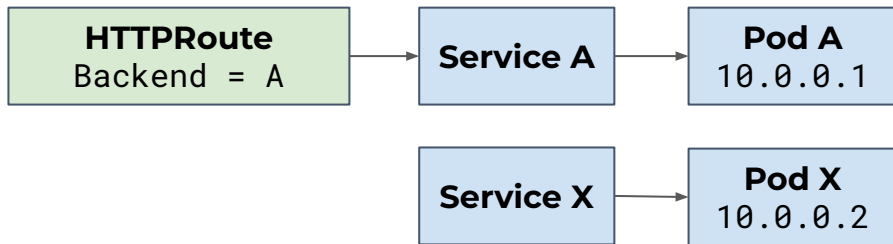


# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes

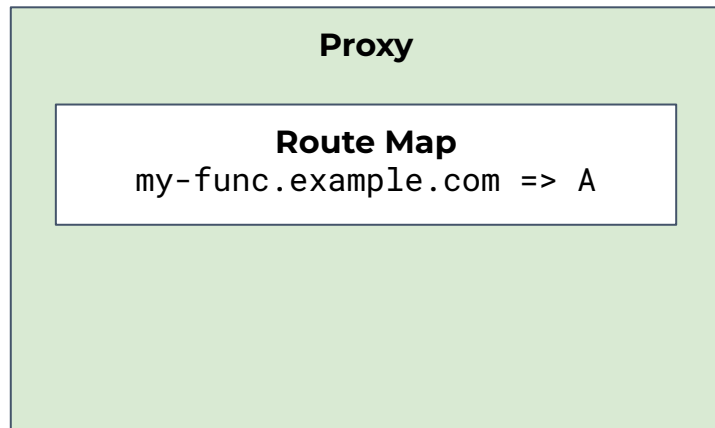
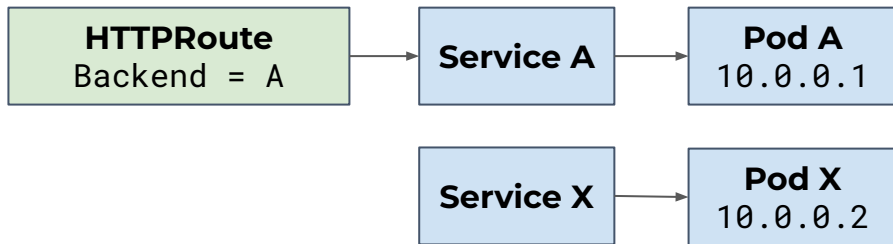
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



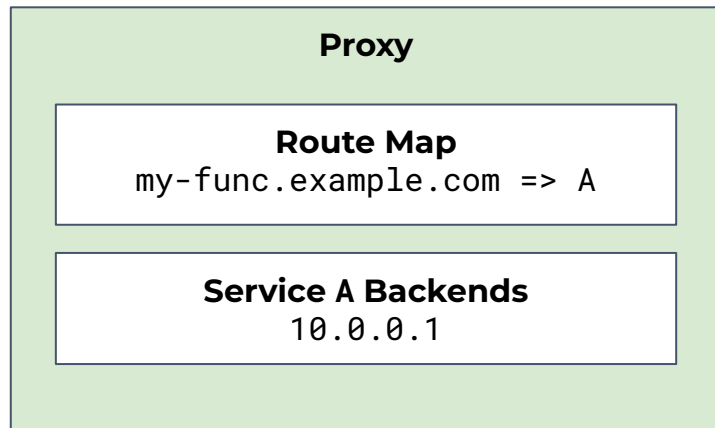
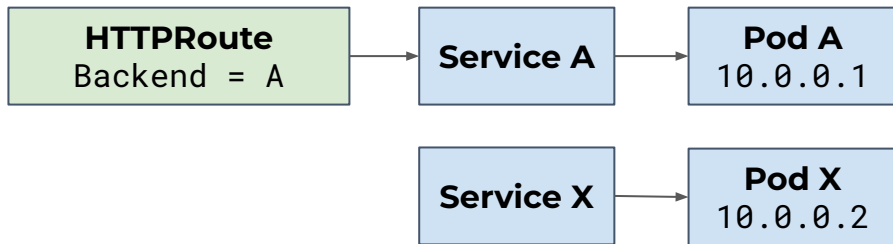
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



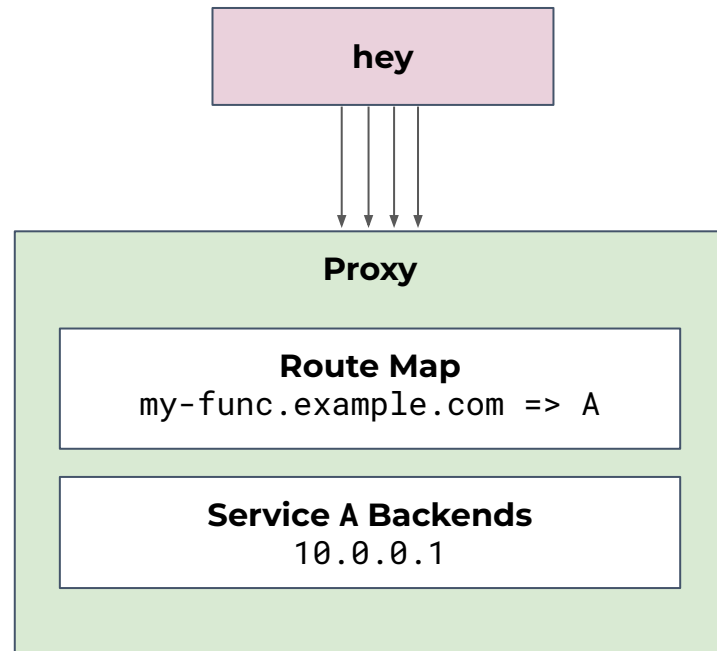
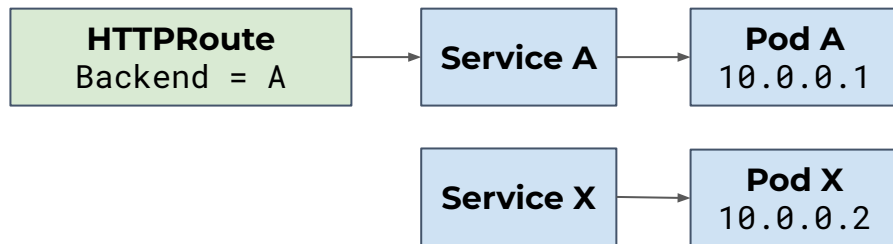
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



# Update Operations are Unsafe

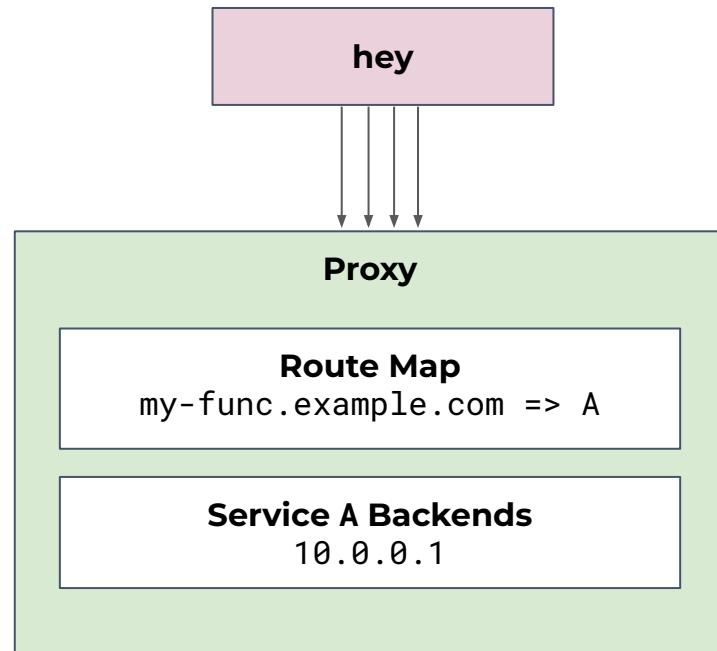
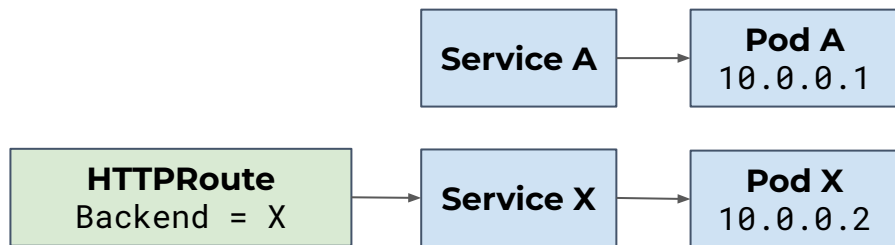
- Ingress doesn't track backend endpoints that aren't reachable via routes





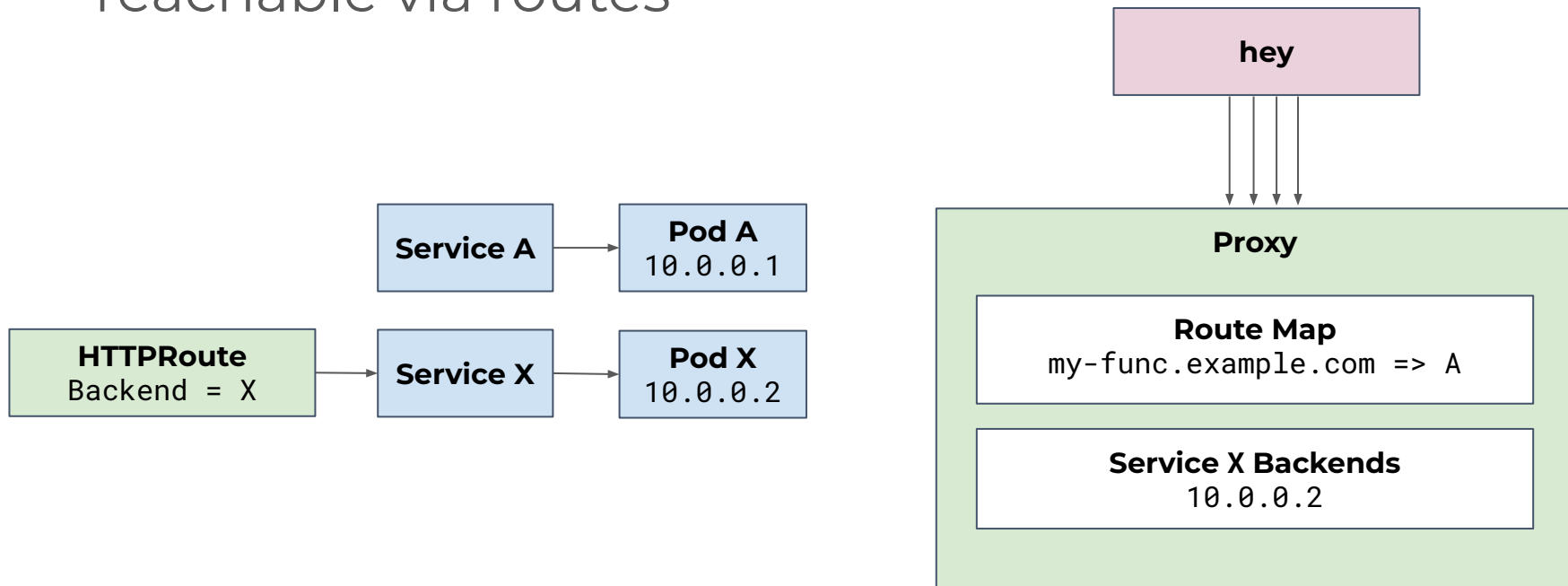
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



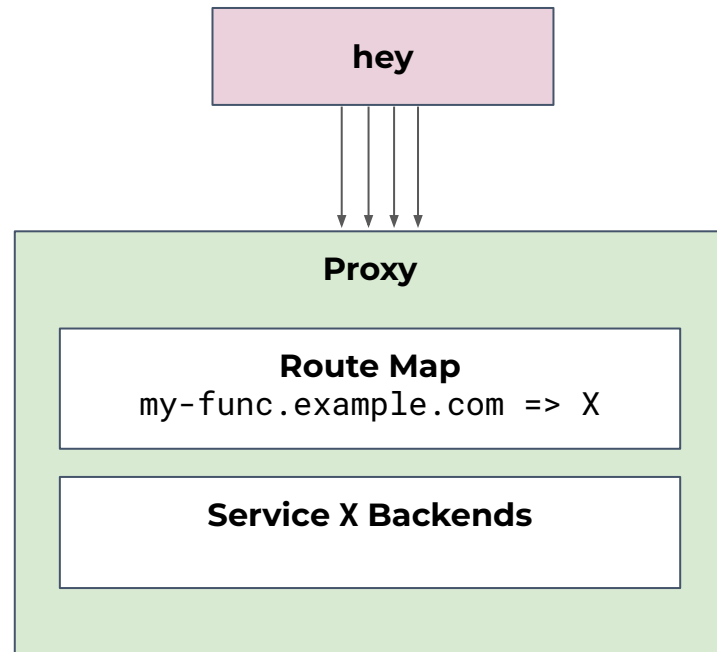
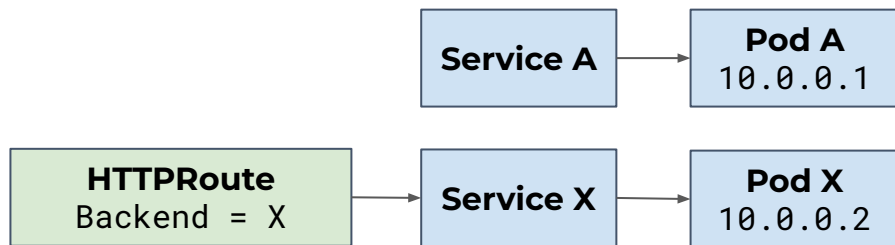
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



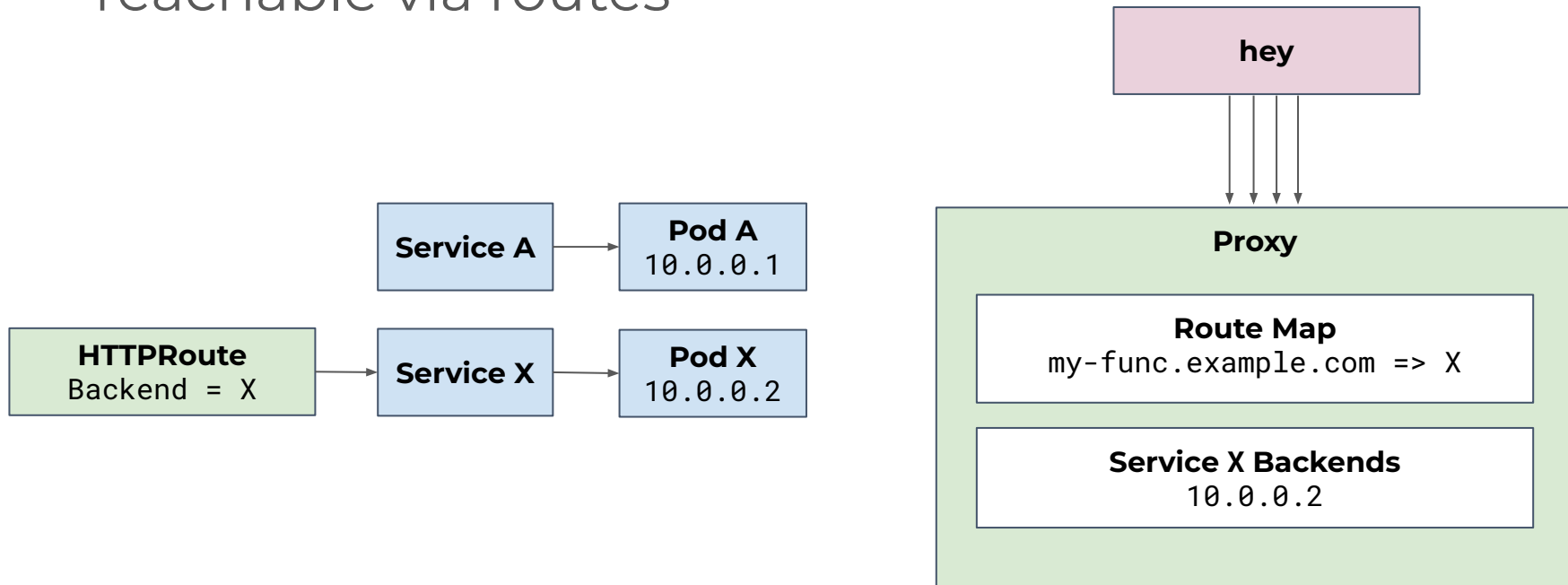
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



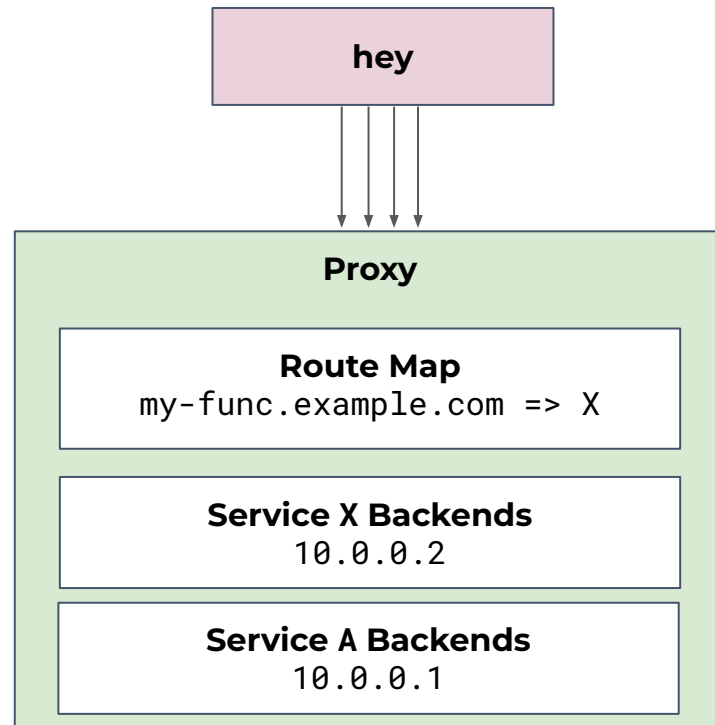
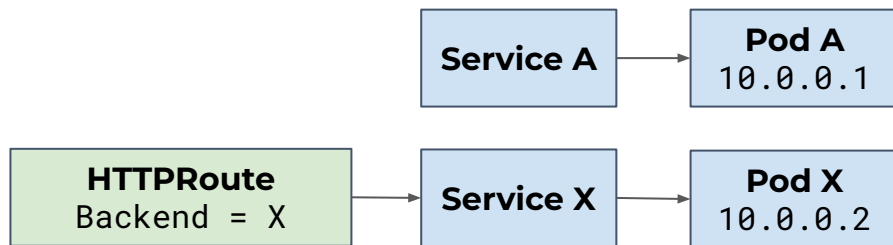
# Update Operations are Unsafe

- Ingress doesn't track backend endpoints that aren't reachable via routes



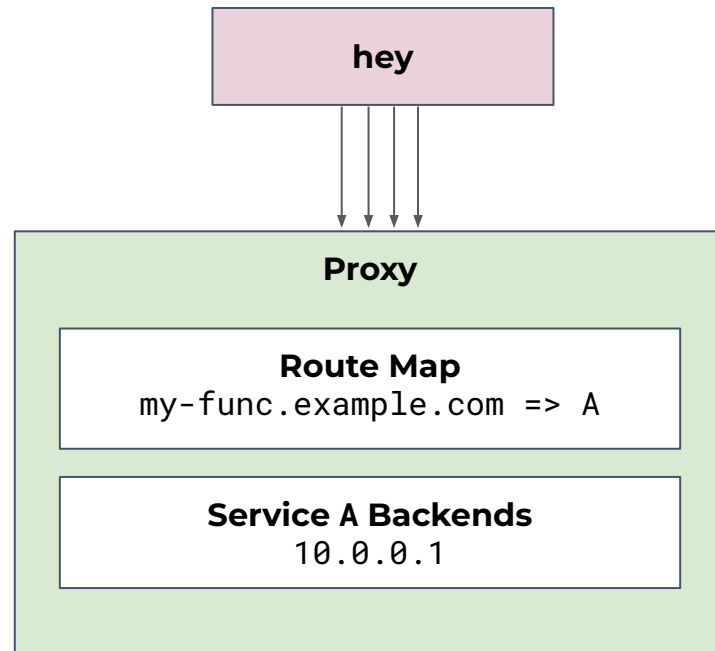
# Update Operations are Unsafe

- Proxy needs to track all potential backends to avoid this



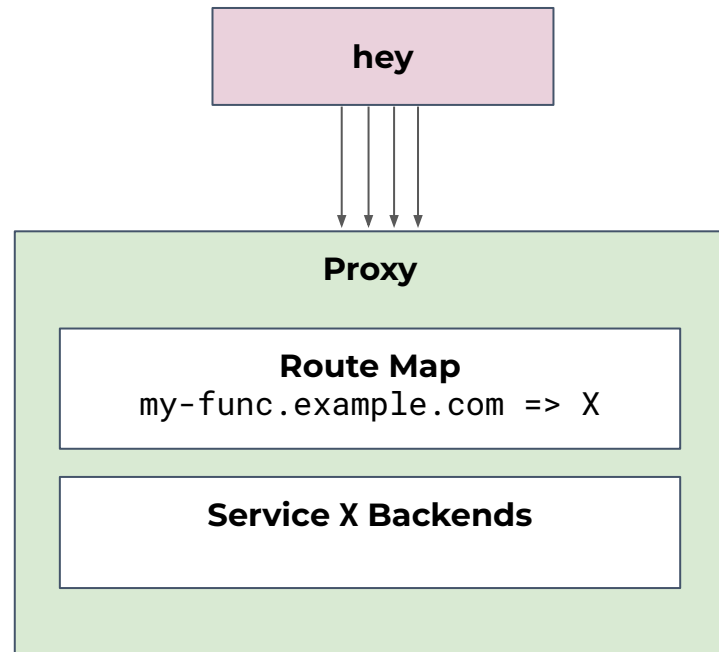
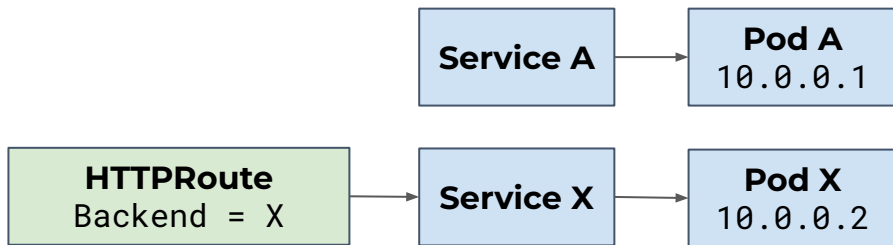
# Update Operations are Unsafe

- Newly deployed healthy backends aren't safe for Routes to reference



# Update Operations are Unsafe

- Newly deployed healthy backends aren't safe for Routes to reference





KubeCon



CloudNativeCon

North America 2024





# shoutout to @mattmoor

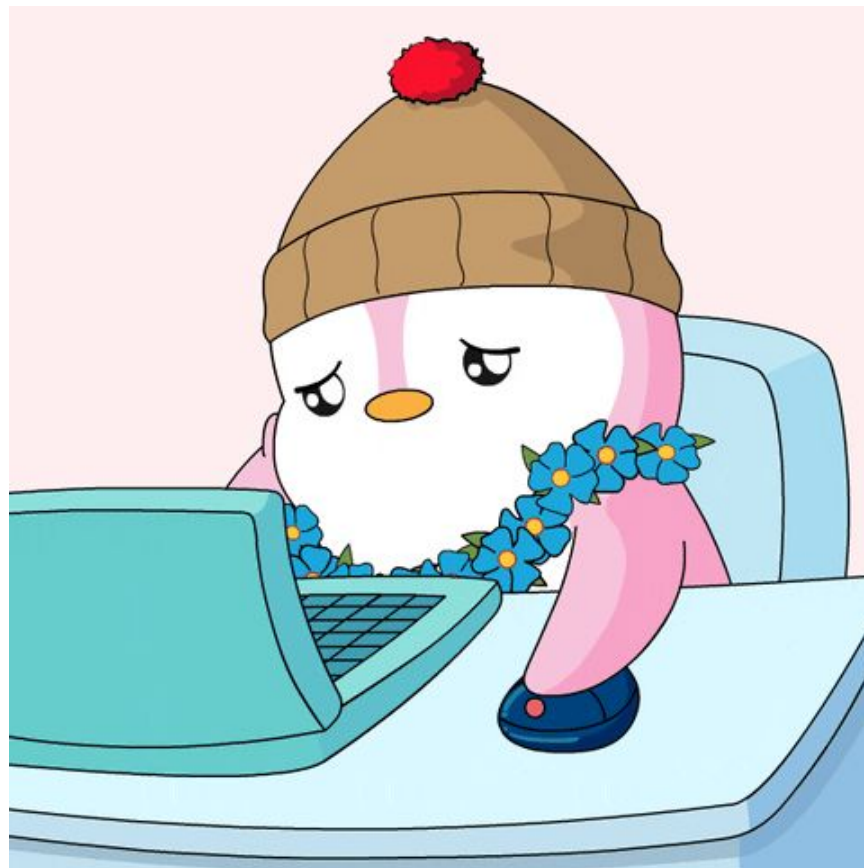


KubeCon

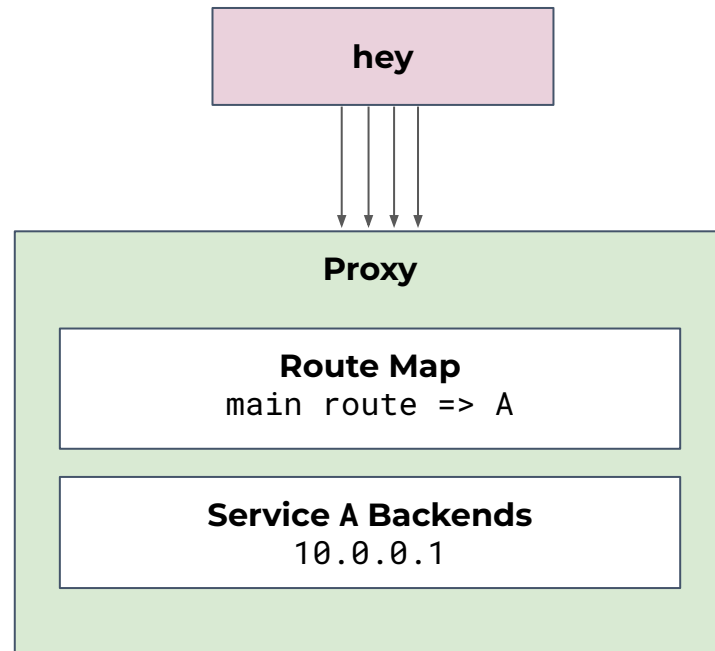
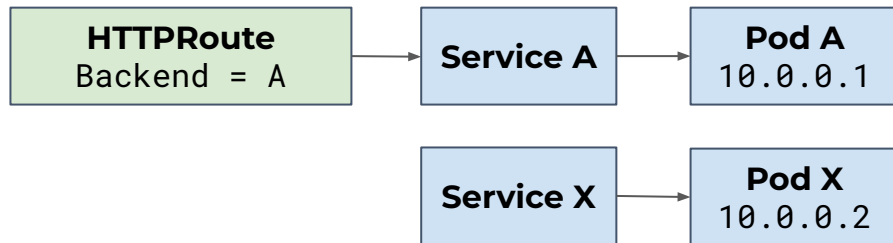


CloudNativeCon

North America 2024

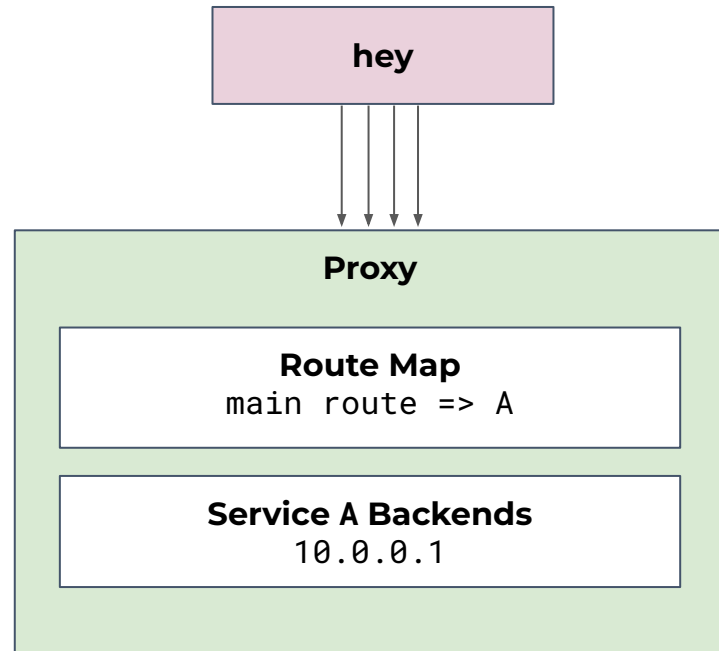
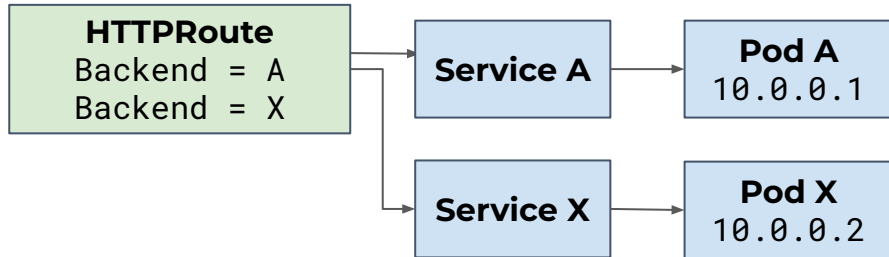


# Three Phase Update



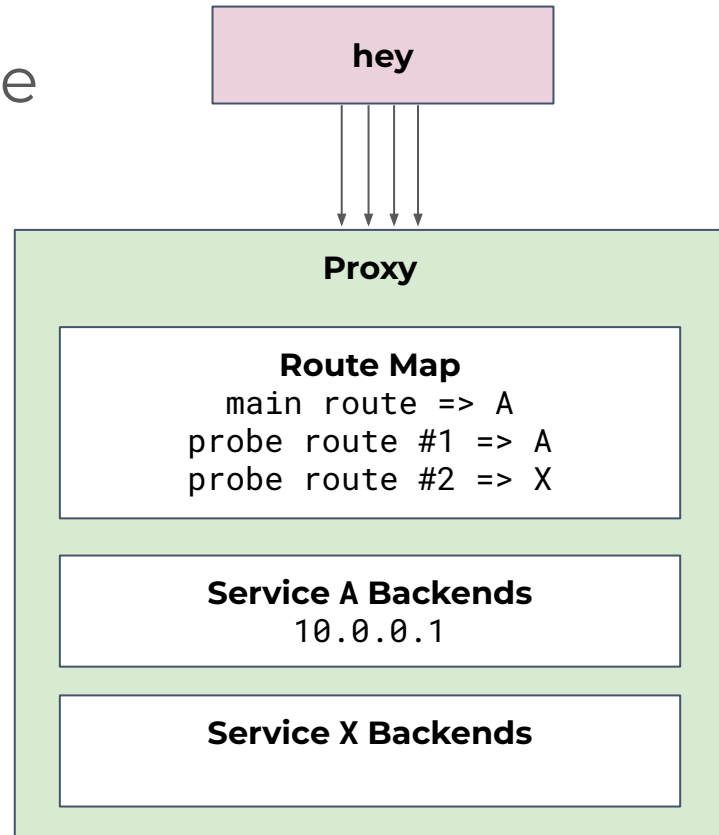
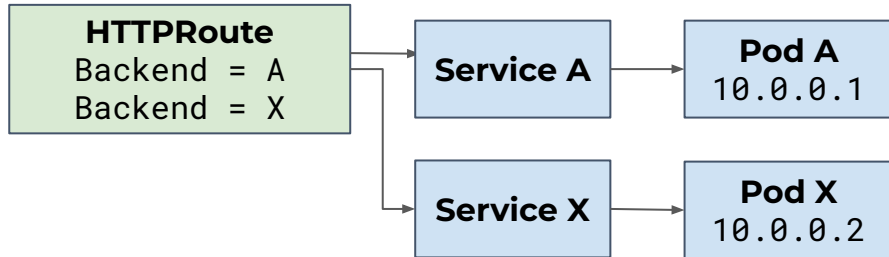
# Three Phase Update

- Add new backends and probe



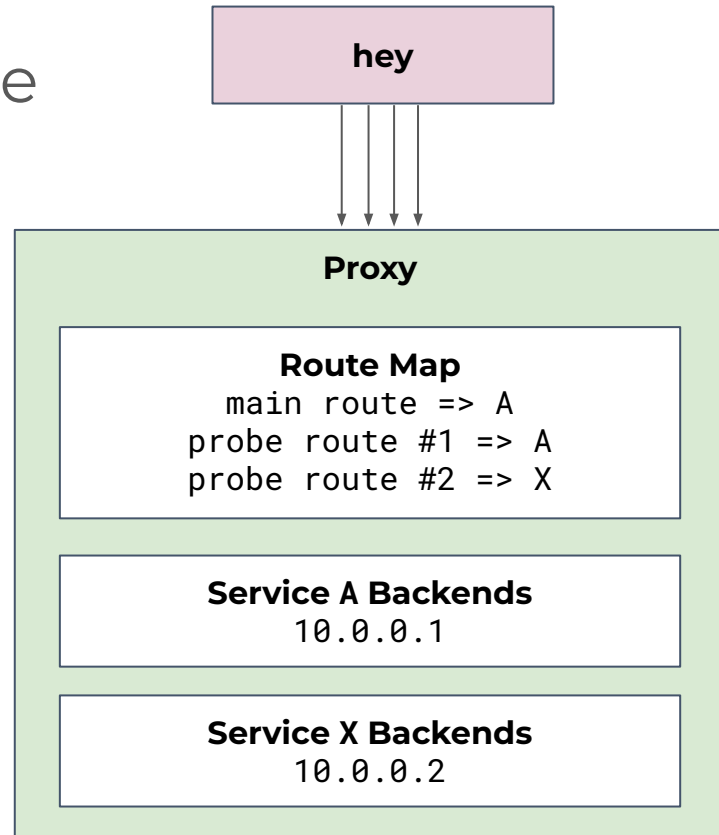
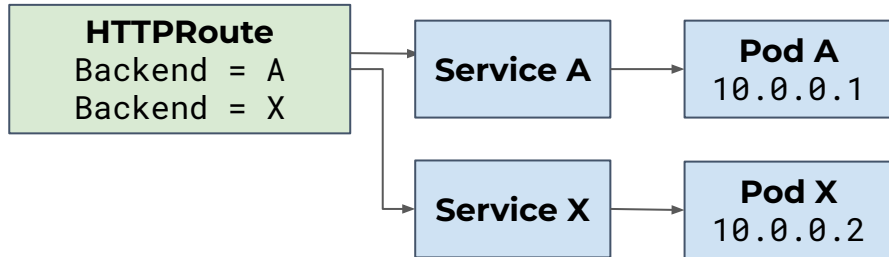
# Three Phase Update

- Add new backends and probe



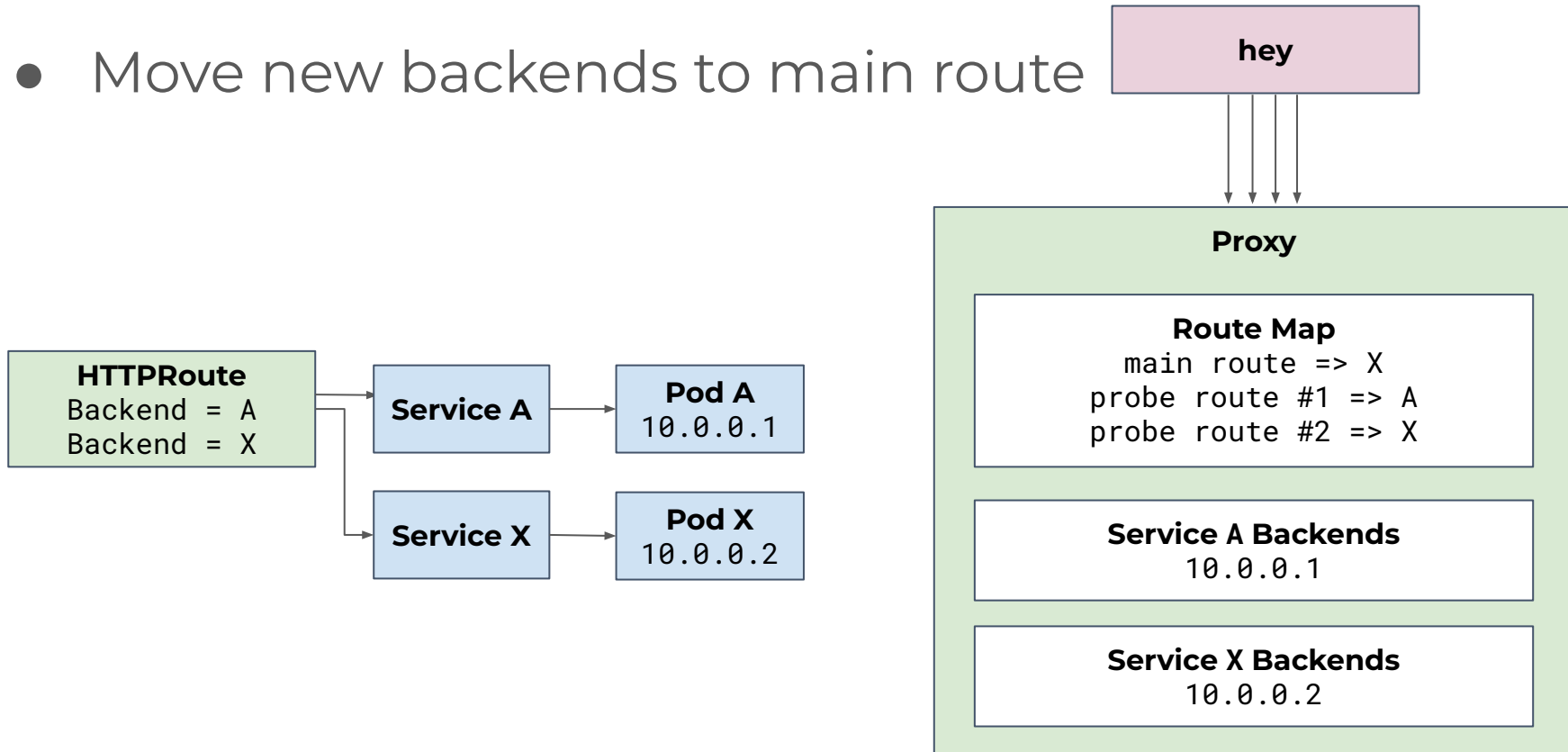
# Three Phase Update

- Add new backends and probe



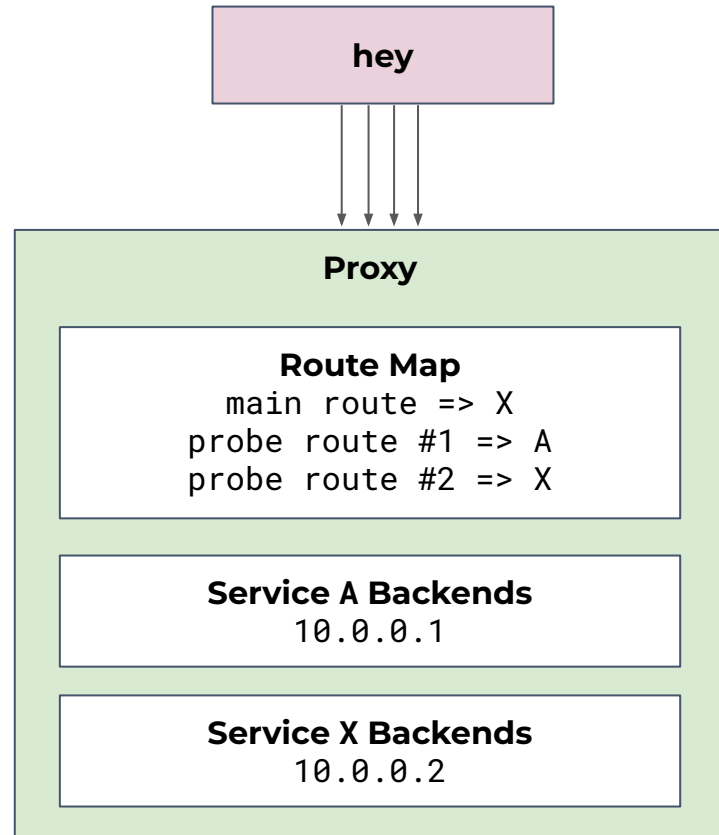
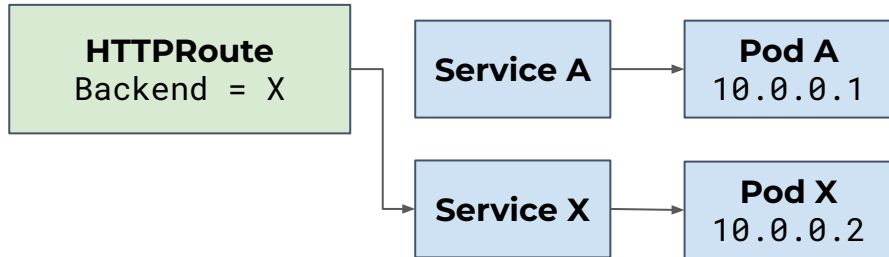
# Three Phase Update

- Move new backends to main route



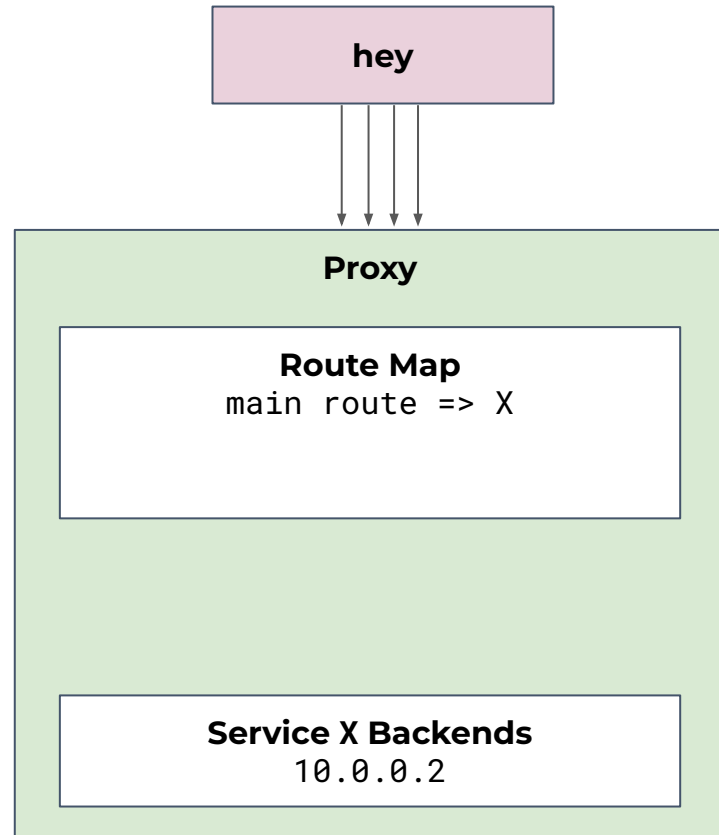
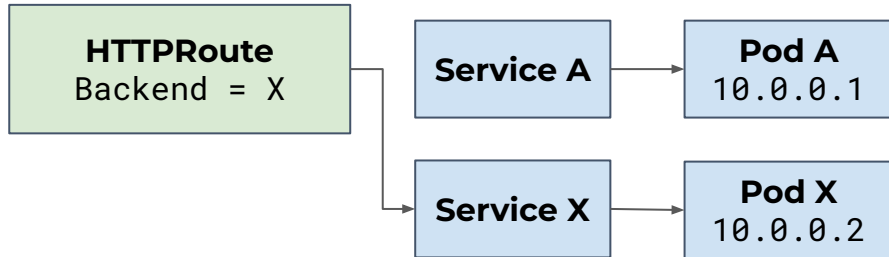
# Three Phase Update

- Drop extra probing routes



# Three Phase Update

- Drop extra probing routes





## Gateway API Endpoint Probing - Design Document

1. Add new backends and probe
2. Move new backends to main route
3. Drop extra probing routes

- Proxies that tract a subset of Services in a cluster are most affected

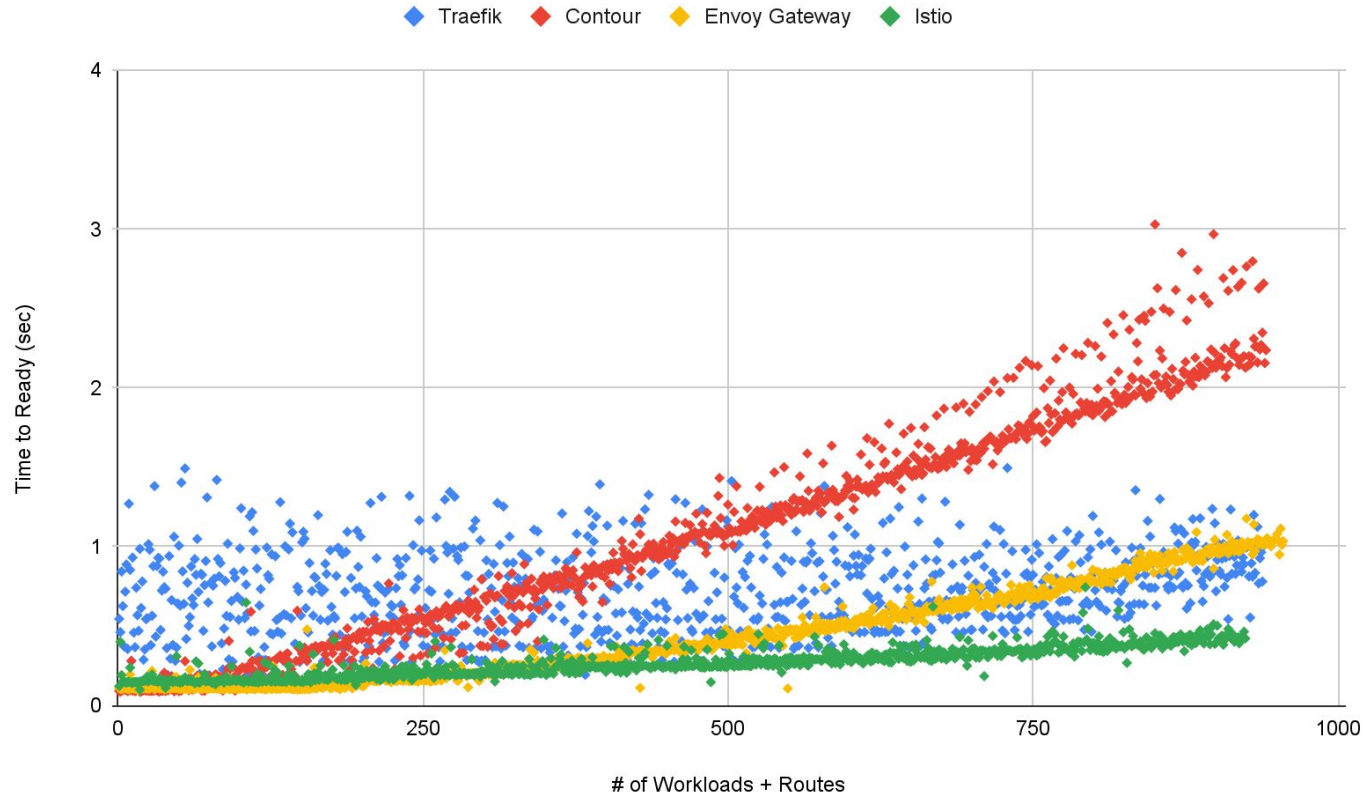
## AND

- If you perform route updates too soon you'll get failures
  - Workload is not Ready
  - Endpoints haven't been propagated to control plane



1. Create a Service & Pod
  - Wait for Endpoints to have a ready address
2. Create an HTTPRoute with a unique domain
3. Poll the route and determine how long it takes to become ready
4. Repeat 1-3 - 1000 times

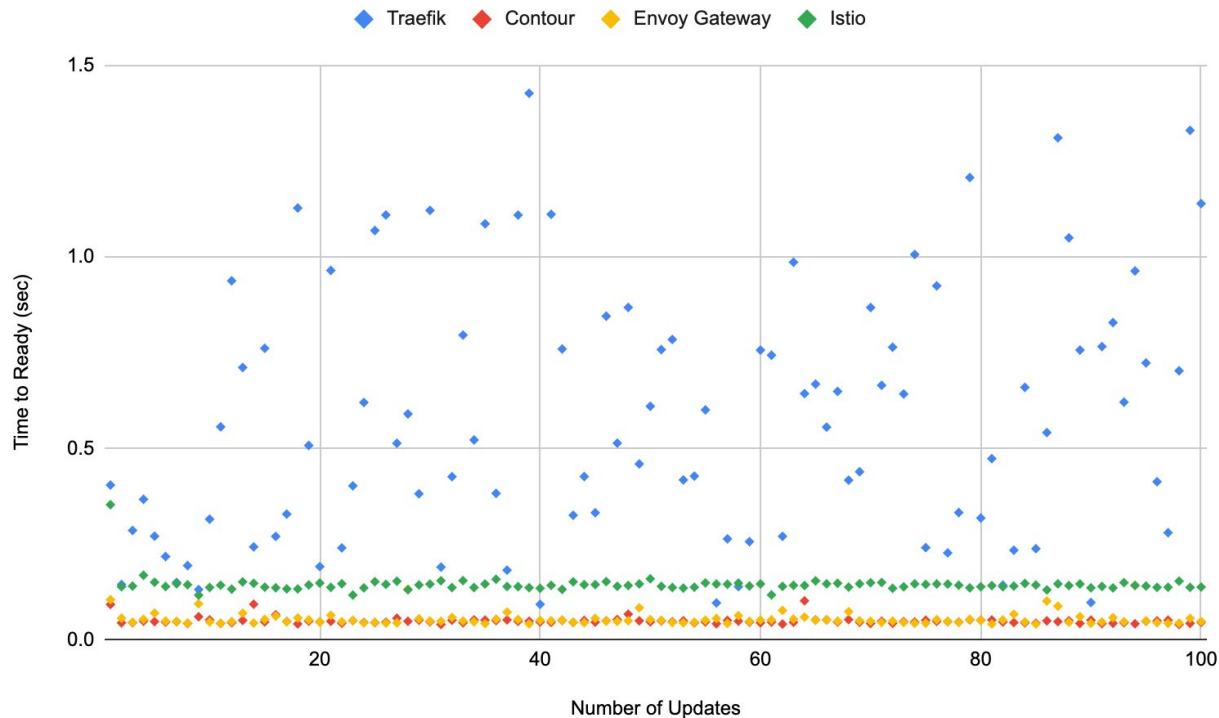
# Readiness Time vs. Creation



# Ready Time - Updating a Single Route

1. Create an initial Service & Pod
2. Create a single HTTPRoute and start polling
3. Create a new Service & Pod
  - Wait for Endpoints to have a ready address
4. Update the HTTPRoute
5. Repeat 4-5 100 times
6. Check poller for dropped traffic

# Readiness Time vs. Updates

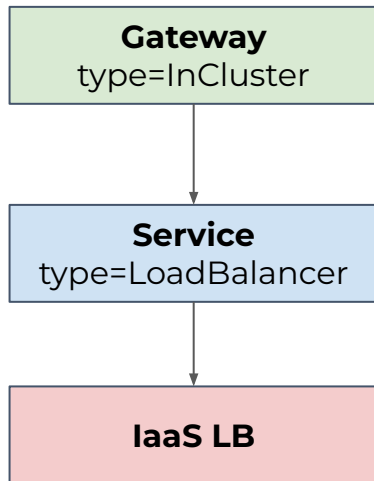


	Traefik	Contour	Envoy Gateway	Istio
SLI	100.00%	99.62%	100.00%	100.00%

- GKE Managed Gateways
  - **gke-17-regional-external-managed**
    - Ready time after CREATE ~ **1m29s**
    - Ready time after UPDATE ~ **1m40s**
      - Continuous Polling - **SLI 64%**
  - **gke-17-gx1b**
    - Similar creation times as above
    - Updates never succeeded under ~ **2min**

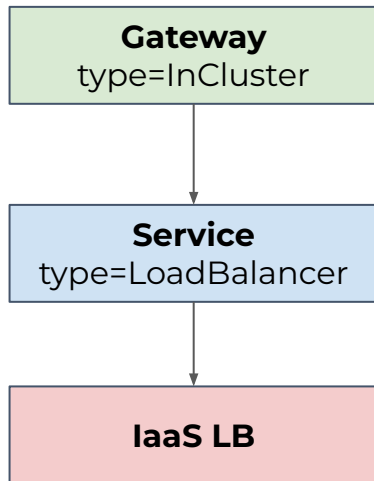


# Underlying K8s Cluster Quirks



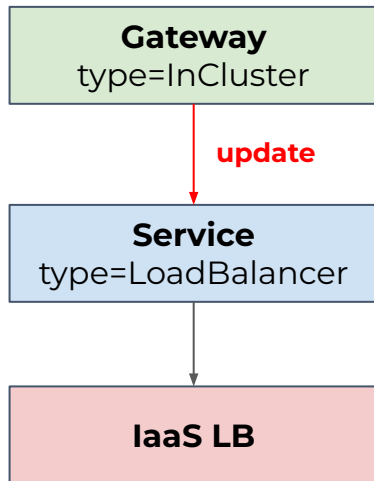
```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
```

# Underlying K8s Cluster Quirks



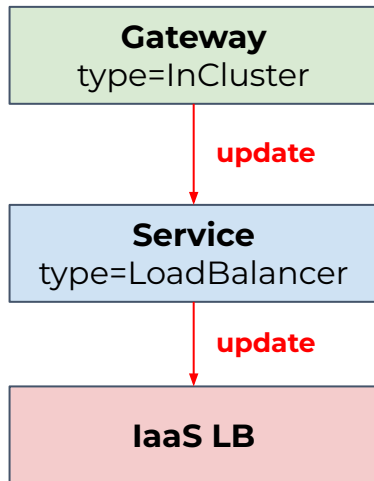
```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTP
      port: 9090
      name: metrics
```

# Underlying K8s Cluster Quirks



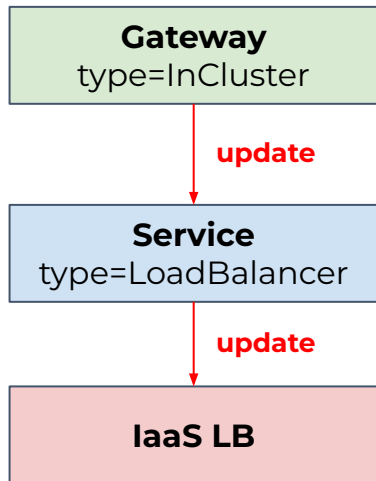
```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTP
      port: 9090
      name: metrics
```

# Underlying K8s Cluster Quirks



```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTP
      port: 9090
      name: metrics
```

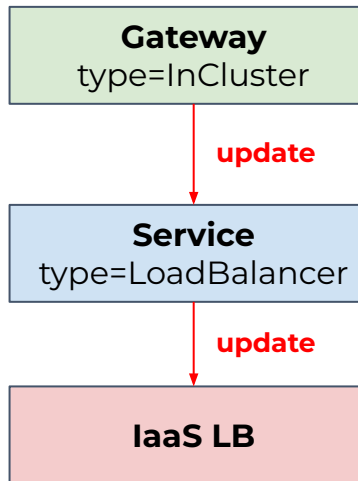
# Underlying K8s Cluster Quirks



```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTP
      port: 9090
      name: metrics
```

- On GKE **ServicePort** changes bring down the whole LB
  - In above example continuous traffic hitting port **80** will fail
  - Upvote issue here **<https://bit.ly/google-fix-my-lb>**

# Underlying K8s Cluster Quirks

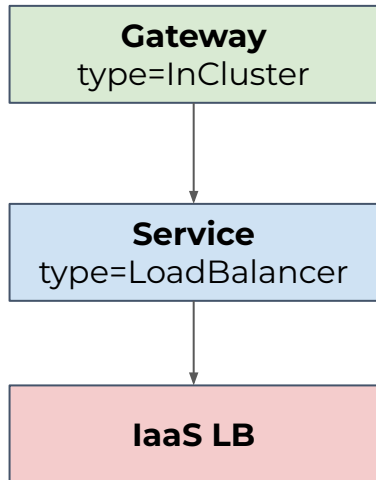


- On GKE Ser
  - In above c
  - Upvote iss



role LB

# Underlying K8s Cluster Quirks



```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTP
      port: 9090
      name: metrics
  infrastructure:
    labels: {}
    annotations: {}
```

<https://gateway-api.sigs.k8s.io/geps/gep-1867/>

# Backend Protocol Selection

## HTTPRoute

```
apiVersion: v1
kind: Service
metadata:
  name: store
spec:
  selector:
    app: store
  ports:
    - protocol: TCP
      appProtocol: kubernetes.io/h2c
      port: 80
```

- [kubernetes.io/h2c](https://kubernetes.io/h2c) - HTTP/2 Prior Knowledge
- [kubernetes.io/ws](https://kubernetes.io/ws) - WebSocket over HTTP

<https://gateway-api.sigs.k8s.io/guides/backend-protocol/>



# GEP-1713: Merging Listeners

```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTPS
      port: 443
      name: foo
      hostname: foo.example.com
      tls:
        certificateRefs:
          - kind: Secret
            group: ""
            name: foo-example-com-cert
```

- Automatic Certificate Provisioning for HTTPS

- Automatic Certificate Provisioning for HTTPS
  - Delegate to cert-manager
  - HTTP-01 is most popular challenge type
    - Doesn't support wildcard certificates

# GEP-1713: Merging Listeners

```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTPS
      port: 443
      name: foo
      hostname: foo.example.com
      tls:
        certificateRefs:
          - kind: Secret
            group: ""
            name: foo-example-com-cert
```

# GEP-1713: Merging Listeners

```
apiVersion: gateway.networking.k8s.io/v1
kind: Gateway
metadata:
  name: prod-web
spec:
  gatewayClassName: example
  listeners:
    - protocol: HTTP
      port: 80
      name: web
    - protocol: HTTPS
      port: 443
      name: foo
      hostname: foo.example.com
      tls:
        certificateRefs:
          - kind: Secret
            group: ""
            name: foo-example-com-cert
    - protocol: HTTPS
      port: 443
      name: bar
      hostname: bar.example.com
      tls:
        certificateRefs:
          - kind: Secret
            group: ""
            name: bar-example-com-cert
```

# GEP-1713: Merging Listeners

```
//  
// Implementations MAY merge separate Gateways onto a single s  
// Addresses if all Listeners across all Gateways are compatib  
//  
// Support: Core  
//  
// +listType=map  
// +listMapKey=name  
// +kubebuilder:validation:MinItems=1  
// +kubebuilder:validation:MaxItems=64  
// +kubebuilder:validation:XValidation:message="tls must not b  
// +kubebuilder:validation:XValidation:message="tls mode must  
// +kubebuilder:validation:XValidation:message="hostname must  
// +kubebuilder:validation:XValidation:message="Listener name  
// +kubebuilder:validation:XValidation:message="Combination of  
Listeners []Listener `json:"listeners"`
```

# GEP-1713: Merging Listeners

The screenshot shows a GitHub pull request page for the repository `kubernetes-sigs/gateway-api`. The pull request is titled `GEP-1713: Standard Mechanism to Merge Gateway Listeners (rev 2) #3213`. It is currently open and has 100 conversations, 31 commits, 3 checks, and 2 files changed. The pull request is from `dprotaso:gep-1713-rev-2` to `kubernetes-sigs:main`. A comment by `dprotaso` on Jul 23, edited, provides details about the pull request. The comment includes the following text:

- This replaces [#1863](#)
- What type of PR is this?
- /kind gep
- What this PR does / why we need it:
- Outlines a mechanism to merge Gateway Listeners
- Which issue(s) this PR fixes:
- Fixes [#1713](#)
- Does this PR introduce a user-facing change?:
- NONE

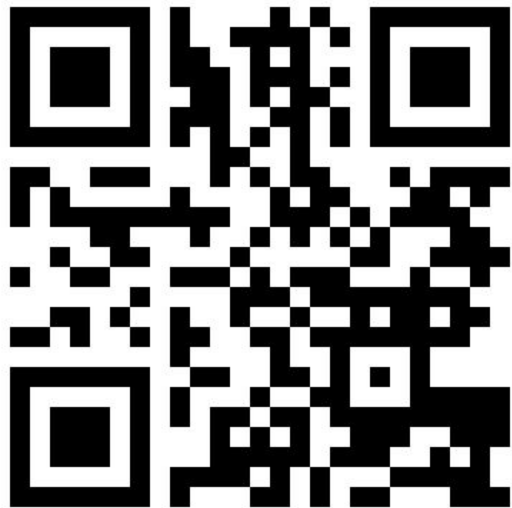
The comment also shows 1 thumbs up and 1 reaction. On the right side of the pull request, there is a list of reviewers: `kate-osborn`, `youngnick`, `guicassolato`, `mikemorris`, `howardjohn`, `robscott`, `LiorLieberman`, and `mlavacca`. Below the reviewers, it says "Still in progress? Convert to draft". Under the "Assignees" section, it says "No one assigned". At the bottom right, there are labels: `cncf-cla: yes`, `do-not-merge/hold`, `priority/important-soon`, and `release-note`.

<https://github.com/kubernetes-sigs/gateway-api/pull/3213>

- Discussions are ongoing!
- Comment on PR/original issue
- Maintainer meeting @ KubeCon
  - Ambassador and Maintainer Lounge: 155D
  - Nov 13 - 3pm



# Questions and Feedback!



[dprotaso](#)

[lintinmybelly](#)

[dprotasowski](#)