

Introduction to Istio

@remkohdev

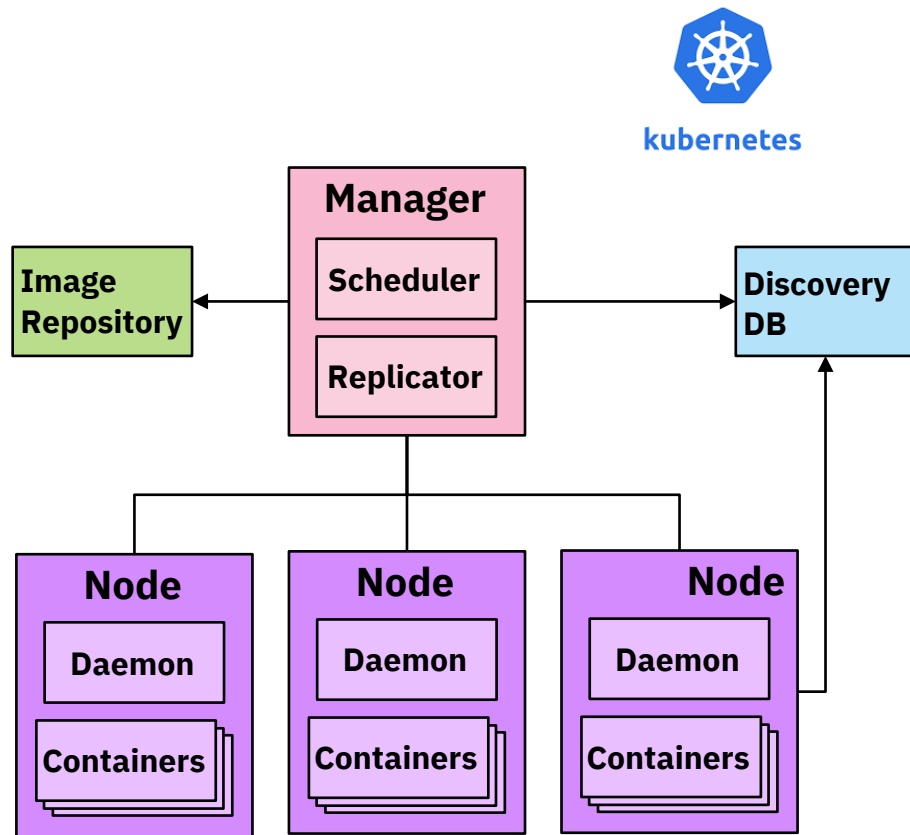
Remko de Knikker

Developer Advocate

Kubernetes

Benefits:

- Automated scheduling and scaling
- Zero downtime deployments
- High availability and fault tolerance
- A/B deployments



Microservices and 12-Factor

- | | | | |
|------|---------------------|-------|-----------------|
| I. | Codebase | VII. | Port Binding |
| II. | Dependencies | VIII. | Concurrency |
| III. | Config | IX. | Disposability |
| IV. | Backing Services | X. | Dev/Prod Parity |
| V. | Build, Release, Run | XI. | Logs |
| VI. | Processes | XII. | Admin processes |

Microservices Challenges

Traffic Management

- How to do canary testing?
- How to do A/B testing?

← All of these problems are common across services no matter the runtime.

Performance, Scalability and Resilience

- How to implement circuit breakers?
- How to test faults in the system?
- How to set rate limits for each service?

← How to solve these problems in a standard way **without changing application code**

Telemetry and Observability

- How to trace a request through my system of multiple services
- How to get centralized logging in a distributed deployment?

Policies and Security

- How to apply policies across services?
- How to manage certificates and trust, and TLS traffic?

Deployment

- How to deploy to multi-cluster environment?
- How to manage complex deployment strategies?

What is a service mesh?

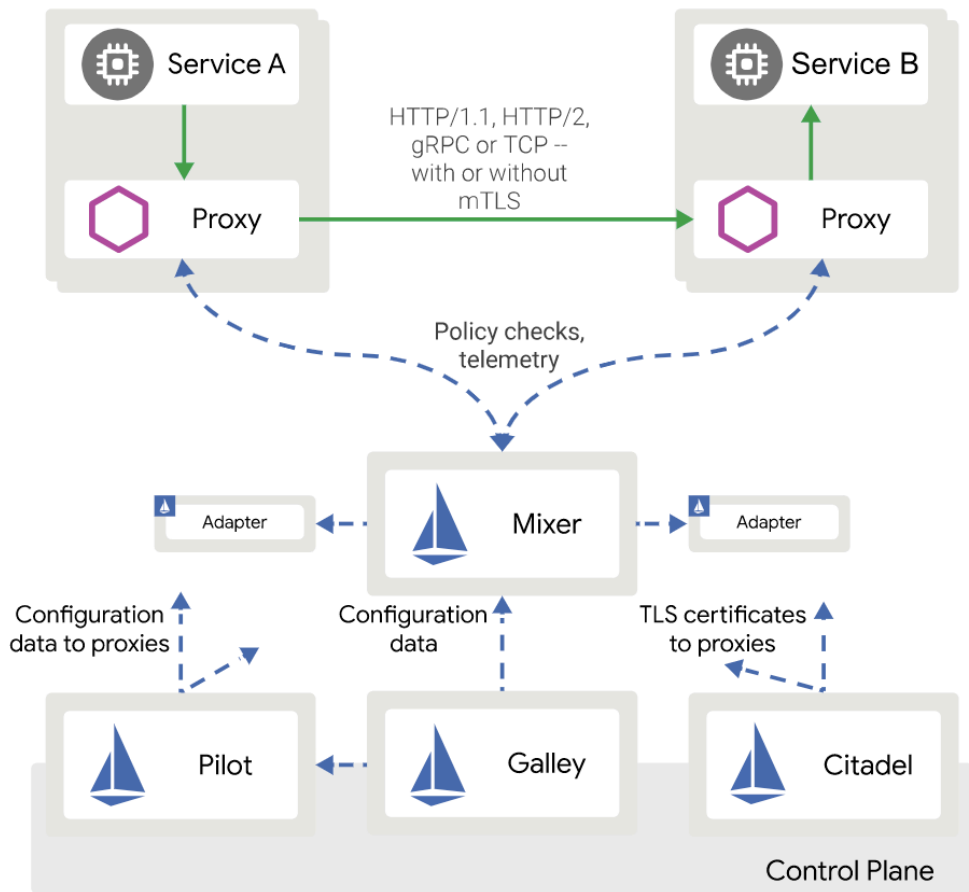
A **service mesh** is a network of microservices and the interactions between them. Unlike other systems for managing this communication, a service mesh is a dedicated infrastructure layer.

Its requirements can include discovery, load balancing, failure recovery, metrics, and monitoring as well as more complex operational requirements, like A/B testing, canary rollouts, rate limiting, access control, and end-to-end authentication.

Istio Architecture

The **data plane** is a set of intelligent proxies (Envoy) deployed as sidecars, which mediate and control all network communication between microservices along with Mixer, a general-purpose policy and telemetry hub.

The **control plane** manages and configures the proxies to route traffic. Additionally, the control plane configures Mixers to enforce policies and collect telemetry.



Istio Architecture

Envoy

- L7 Proxy sidecar,
- Mediates all inbound and outbound traffic,

Mixer

- Enforces access control and usage policies,
- Collects telemetry data from the Envoy proxy sidecar on each request.
- Platform-independent, plug-in model, to interface with a variety of host environments and infrastructure backends.

Pilot

- “Drives” the service mesh by providing service discovery across platforms.
- Traffic management capabilities for intelligent routing (A/B tests, canary rollouts).
- Resilience (timeouts, retries, circuit-breakers and fault injection).
- Coverts that to low-level routing rules, and distributes those to sidecar containers at runtime.

Citadel

- Service-to-service and end-user authentication.
- Encrypt traffic.
- Key and certificate management, distributed as Kubernetes Secrets.

Galley

- Configuration validation, ingestion, processing and distribution, from underlying Kubernetes.

Envoy

Istio uses an extended version of the Envoy proxy, <https://www.envoyproxy.io/>. Envoy is an L7 proxy for the application layer, layer 7 or L7 of the Open Systems Interconnection (OSI) model.

Built-in Envoy features:

- Dynamic service discovery
- Load balancing
- TLS termination
- HTTP/2 and gRPC proxies
- Circuit breakers
- Health checks
- Staged rollouts with %-based traffic split
- Fault injection
- Rich metrics

Mixer

The Envoy sidecar calls Mixer before each request to perform precondition checks, and after each request to report telemetry. One of its key functions is to abstract away the details of different backend systems.

Adapters are individual plug-ins as Go packages.

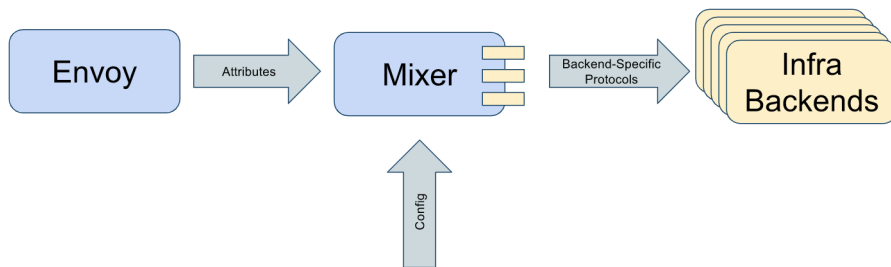
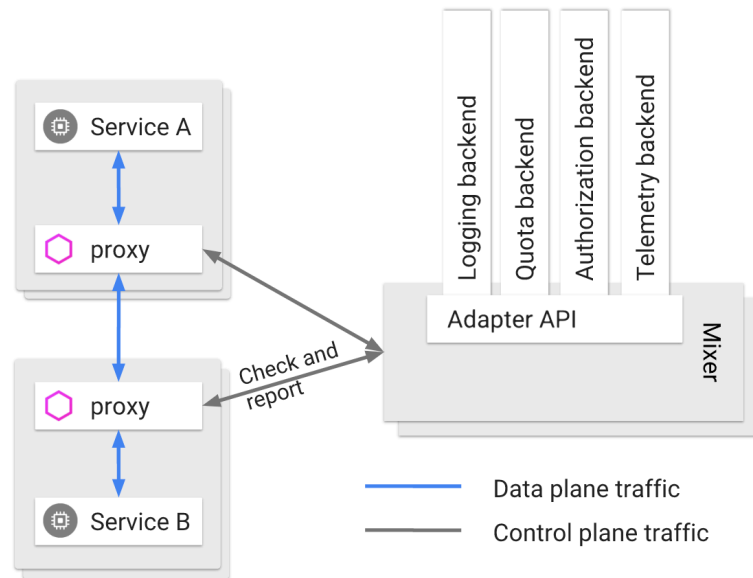
Templates are used to send data to individual adapters.

Instances control which data is delivered to individual adapters.

Rules tell Mixer which instances to send to which handler and when.

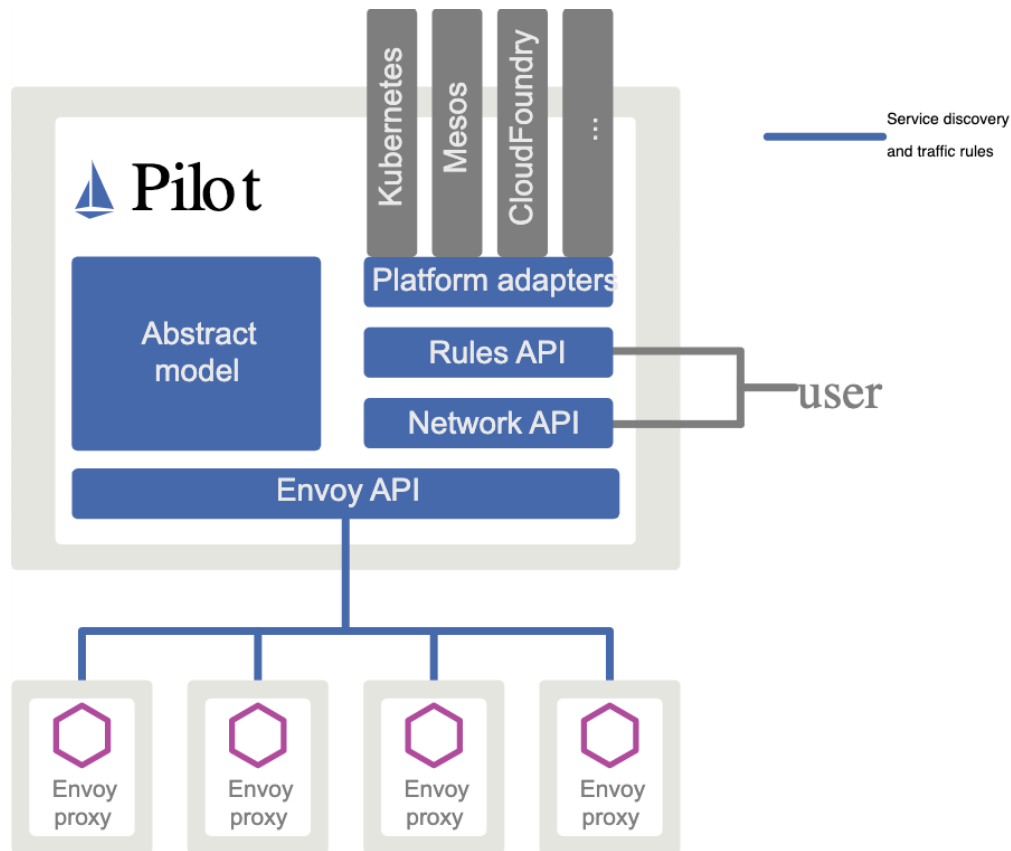
Mixer adapters:

- App Identity and Access
- Datadog
- Fluentd
- Kubernetes env
- Memory quota
- Redis quota
- Stdio
- Zipkin
- Prometheus
- Heapster
- etc



Pilot

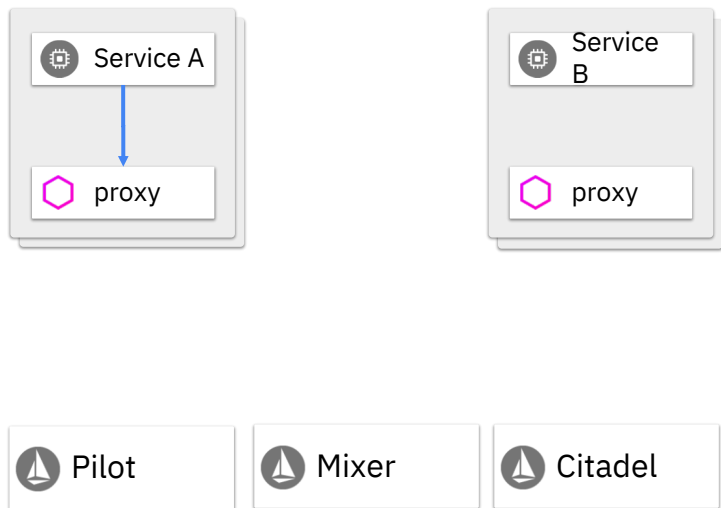
Pilot maintains an **abstract model** of all the services in the mesh, to generate appropriate Envoy-specific configurations to let Envoy proxies know about one another in the mesh through the Envoy API. Use the Traffic Management API to instruct Pilot to refine the Envoy configuration.



Lifecycle of a request



Lifecycle of a request

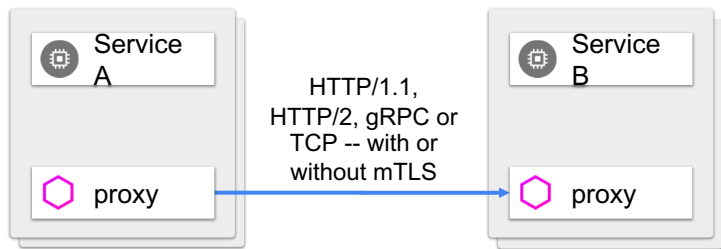


Service A places a call to <http://service-b>.

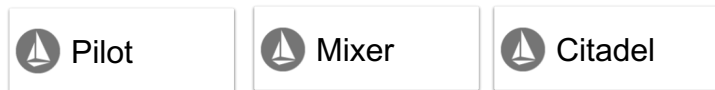
Client-side Envoy intercepts the call.

Envoy consults config (previously received from Pilot) to know how/where to route call to service B (taking into account service discovery, load balancing, and routing rules), forwards the call to the right destination.

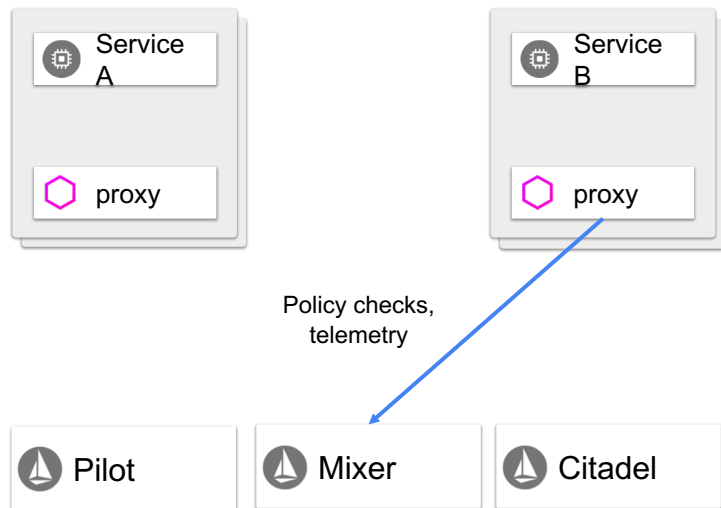
Lifecycle of a request



Envoy forwards request to appropriate instance of service B. There, the Envoy proxy deployed with the service intercepts the call.

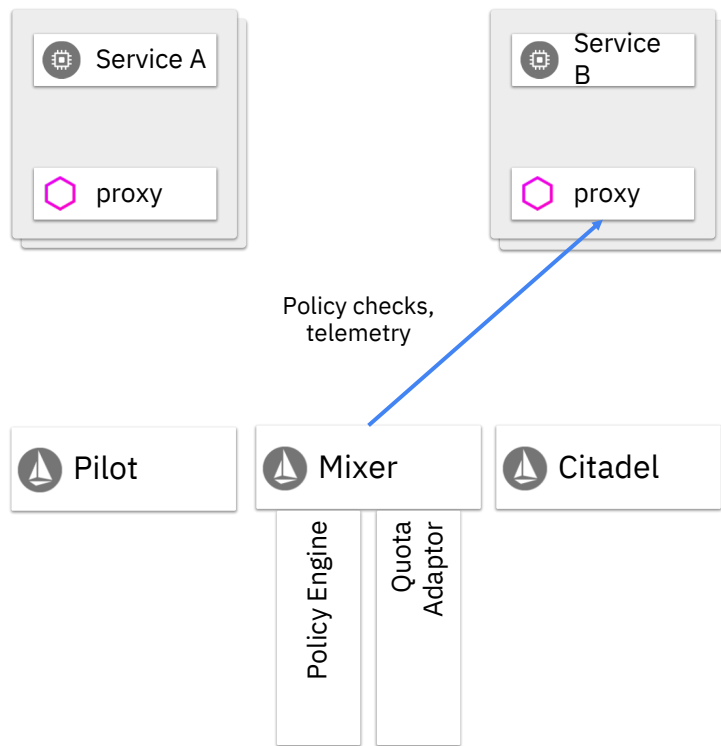


Lifecycle of a request



Server-side Envoy checks with Mixer to validate that call should be allowed (ACL check, quota check, etc).

Lifecycle of a request

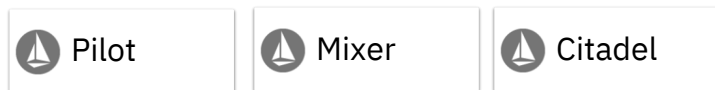


Mixer checks with appropriate adaptors (policy engine, quota adaptor) to verify that the call can proceed and returns true/false to Envoy

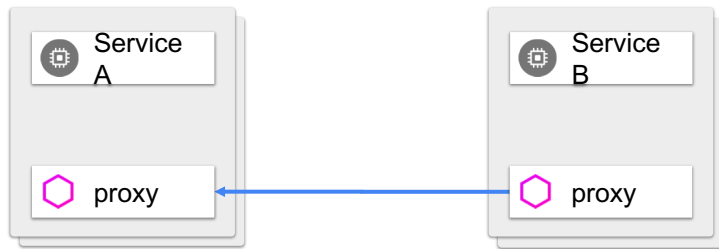
Lifecycle of a request



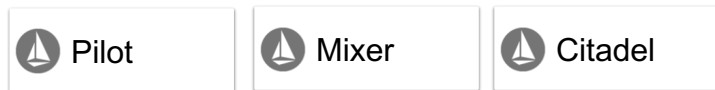
Server-side Envoy forwards request to service B, which process request and returns response



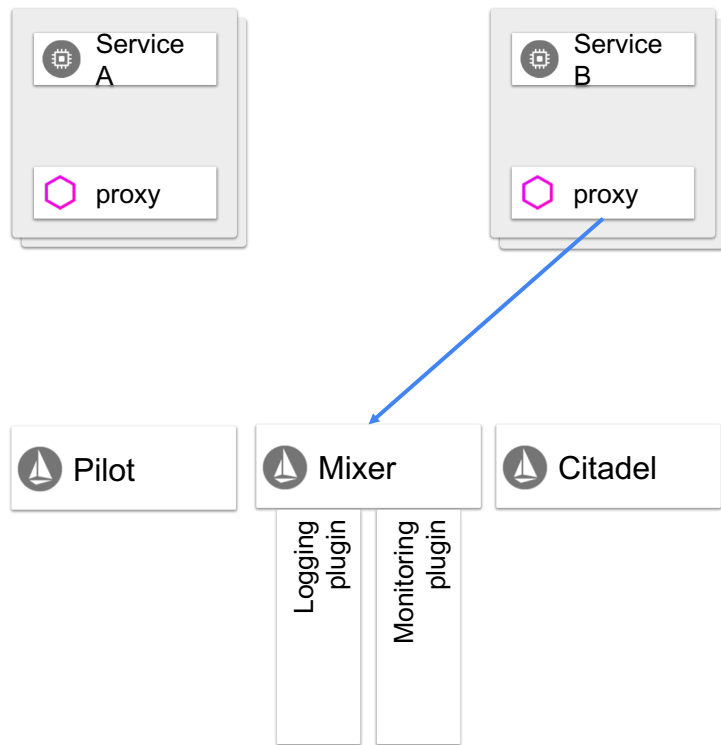
Lifecycle of a request



Envoy forwards response to the original caller, where response is intercepted by Envoy on the caller side.



Lifecycle of a request

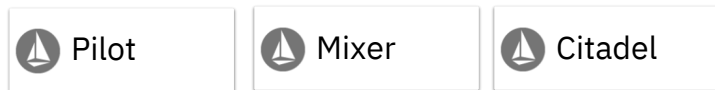


Envoy reports telemetry to Mixer, which in turn notifies appropriate plugins

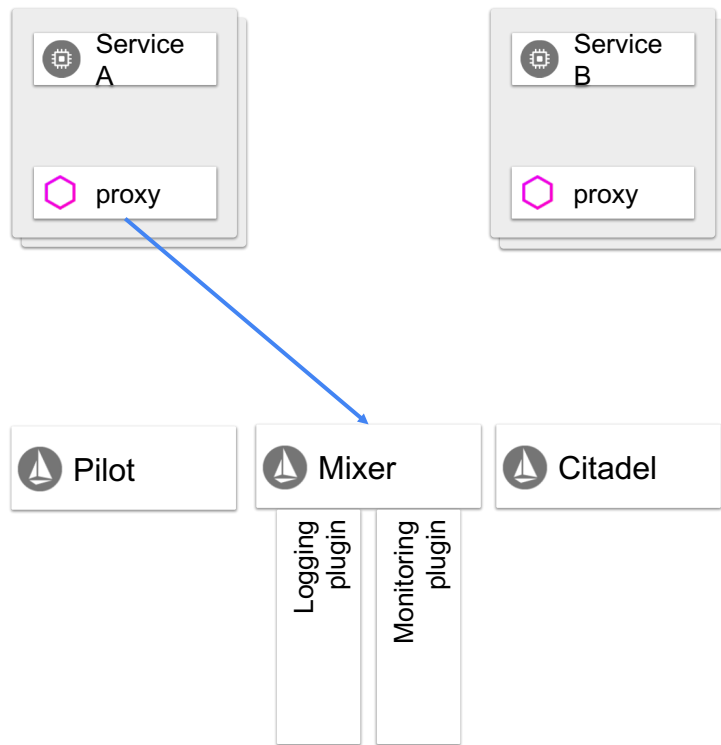
Lifecycle of a request



Client-side Envoy forwards response to original caller.

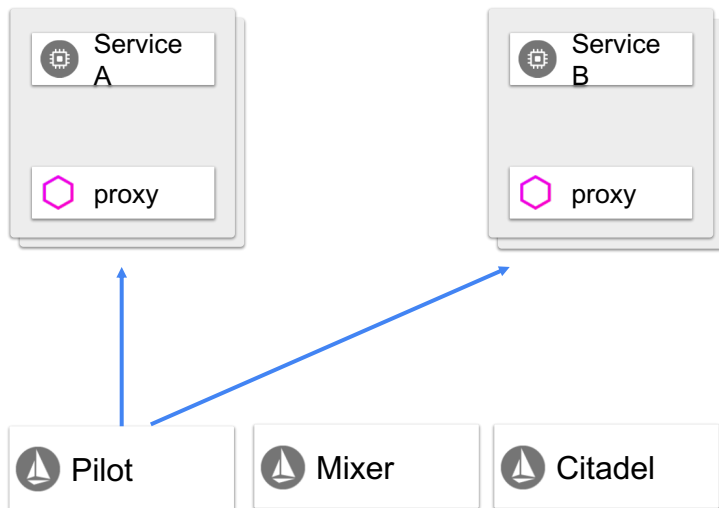


Lifecycle of a request



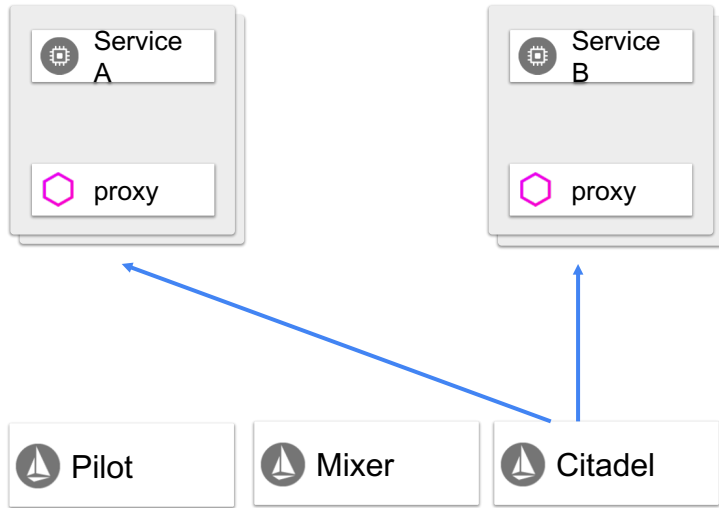
Client-side Envoy reports telemetry to Mixer (including client-perceived latency), which in turn notifies appropriate plugins

At Anytime



Pilot listens to your configuration, such as routing rules, circuit breaking and fault injection. Converts to low-level config (routing rules) and distributes to proxy side cars

At Anytime

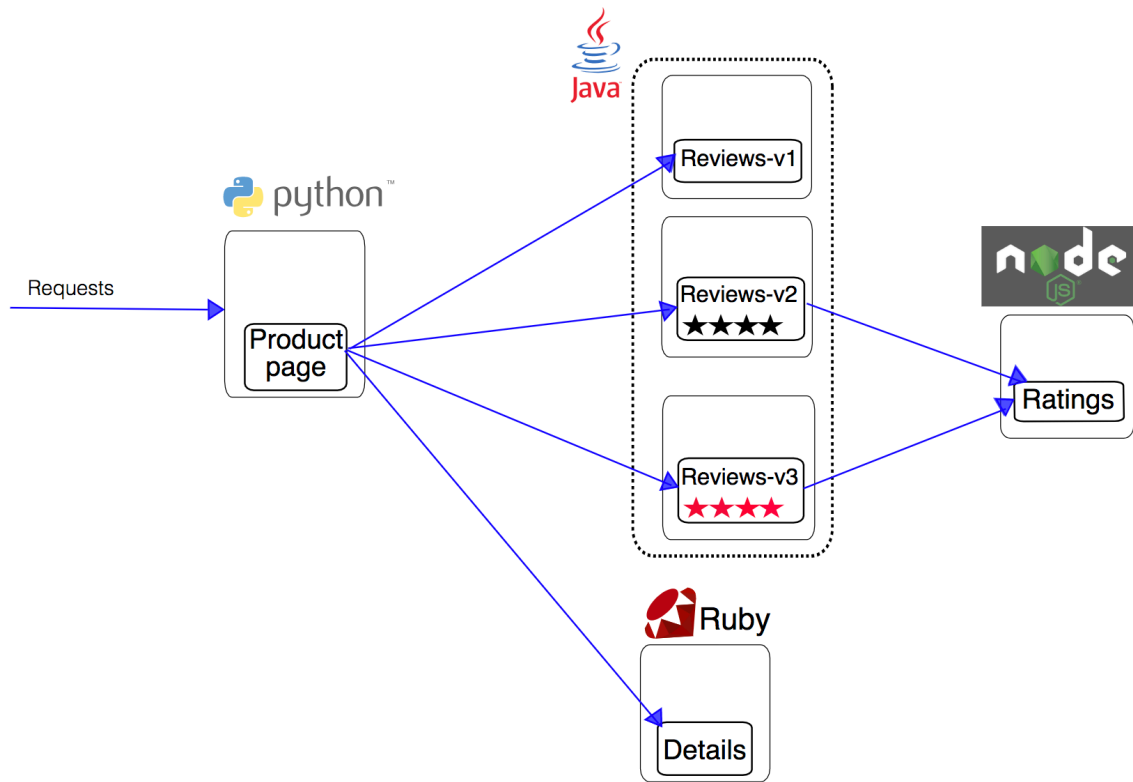


Citadel distributes keys and certificates as Kubernetes Secrets available to sidecar containers

Envoy Proxy Sidecar

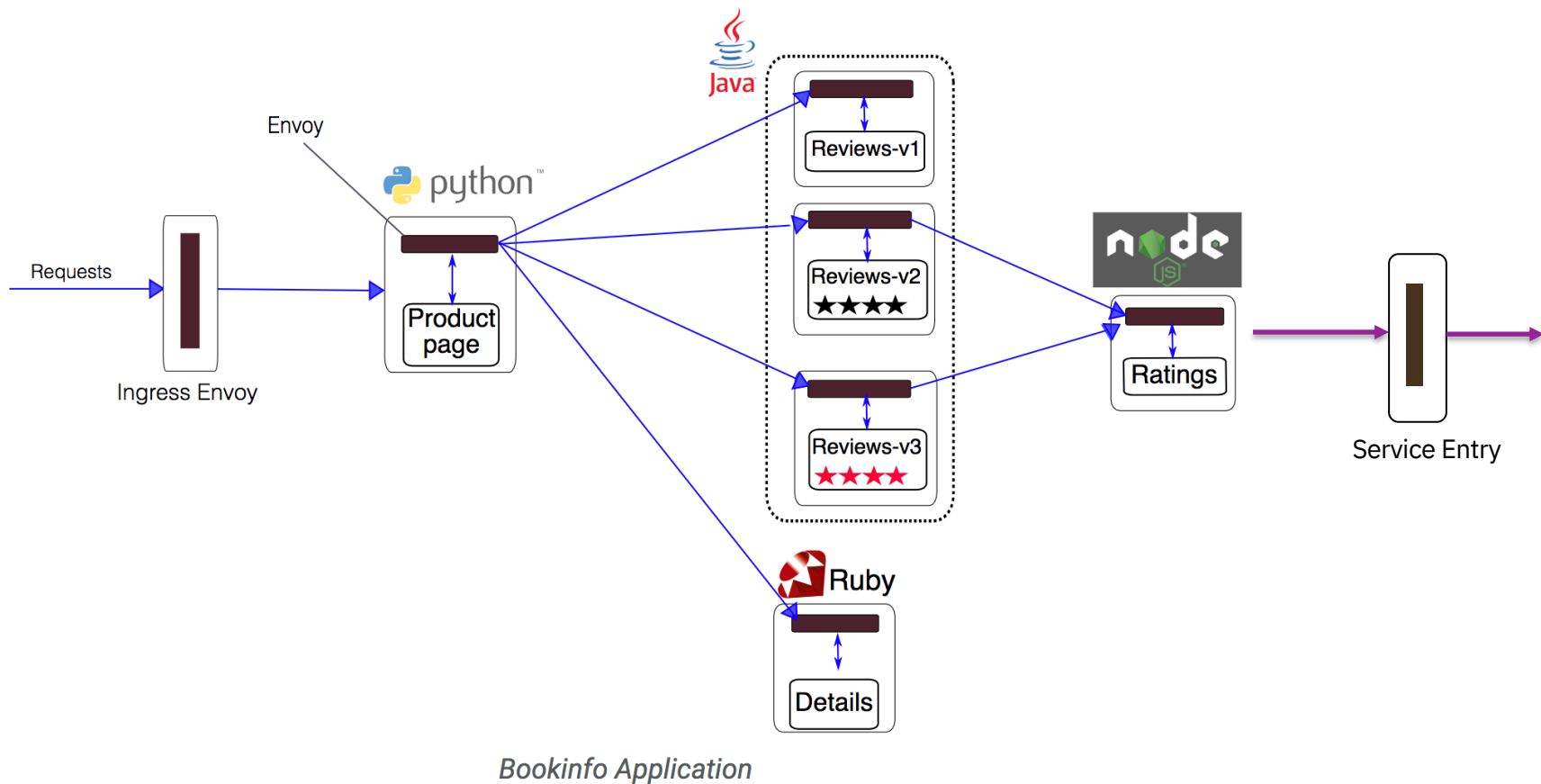
- To create a service mesh with Istio, you update the deployment of the pods to add the Istio Proxy (based on the Lyft Envoy Proxy) as a side car to each pod. The Proxy is then run as a separate container that manages all communication with that pod.
 - Manual Approach
 - Automated Approach
- Istio modifies the individual pods and not deployments

Bookinfo Sample Application (without Istio)




Bookinfo Application without Istio

Bookinfo Sample Application (with Istio)



Request Routing

```
1  apiVersion: networking.istio.io/v1alpha3
2  kind: VirtualService
3  metadata:
4    name: reviews
5  spec:
6    hosts:
7      - reviews
8    http:
9      - match:
10        - headers:
11          end-user:
12            exact: jason
13          route:
14            - destination:
15              host: reviews
16              subset: v2
17        - route:
18          - destination:
19            host: reviews
20            subset: v1
```



Canary Testing

Route user:jason to reviews:v2
Others still get reviews:v1

Traffic Shifting

```
1  apiVersion: networking.istio.io/v1alpha3
2  kind: VirtualService
3  metadata:
4    name: reviews
5  spec:
6    hosts:
7      - reviews
8    http:
9      - route:
10        - destination:
11          host: reviews
12          subset: v1
13          weight: 50
14        - destination:
15          host: reviews
16          subset: v3
17          weight: 50
```

50% -> v1

50% -> v3

Rate Limits

```
1  apiVersion: "config.istio.io/v1alpha2"
2  kind: memquota
3  metadata:
4    name: handler
5    namespace: istio-system
6  spec:
7    quotas:
8    - name: requestcount.quota.istio-system
9      maxAmount: 5000
10     validDuration: 1s
11     # The first matching override is applied.
12     # A requestcount instance is checked against override dimensions.
13     overrides:
14     # The following override applies to 'ratings' when
15     # the source is 'reviews'.
16     - dimensions:
17       destination: ratings
18       source: reviews
19       maxAmount: 1
20       validDuration: 1s
21     # The following override applies to 'ratings' regardless
22     # of the source.
23     - dimensions:
24       destination: ratings
25       maxAmount: 100
```

5000 requests per 1s
ratings: 100 requests per 1s

Circuit Breaking


```
1  apiVersion: networking.istio.io/v1alpha3
2  kind: DestinationRule
3  metadata:
4    name: httpbin
5  spec:
6    host: httpbin
7    trafficPolicy:
8      connectionPool:
9        tcp:
10         maxConnections: 1
11        http:
12         http1MaxPendingRequests: 1
13         maxRequestsPerConnection: 1
14      outlierDetection:
15        consecutiveErrors: 1
16        interval: 1s
17        baseEjectionTime: 3m
18        maxEjectionPercent: 100
```

Max 1 concurrent
connection & request

Delay Injection


Inject 7 second delay

```
1  apiVersion: networking.istio.io/v1alpha3
2  kind: VirtualService
3  metadata:
4    name: ratings
5  spec:
6    hosts:
7    - ratings
8    http:
9    - match:
10      - headers:
11        end-user:
12          exact: jason
13        fault:
14          delay:
15            percent: 100
16            fixedDelay: 7s
17        route:
18          - destination:
19              host: ratings
20              subset: v1
21      - route:
22        - destination:
23            host: ratings
24            subset: v1
```



Fault Injection

```
1  apiVersion: networking.istio.io/v1alpha3
2  kind: VirtualService
3  metadata:
4    name: ratings
5  spec:
6    hosts:
7    - ratings
8    http:
9    - match:
10      - headers:
11        end-user:
12          exact: jason
13        fault:
14          abort:
15            percent: 100
16            httpStatus: 500
17        route:
18        - destination:
19            host: ratings
20            subset: v1
21      - route:
22        - destination:
23            host: ratings
24            subset: v1
```



jason: Return with Error 500

Istio Gateway

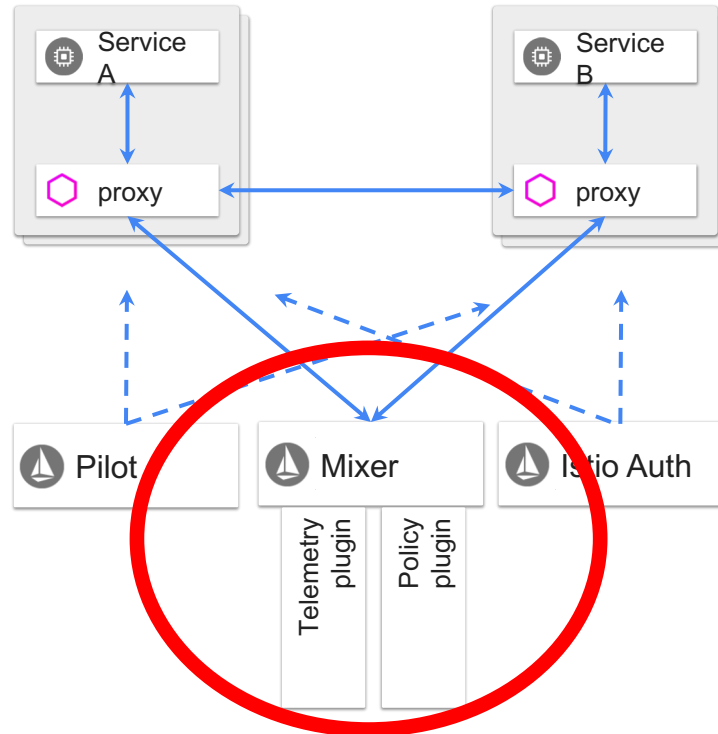


Traffic Management

Telemetry

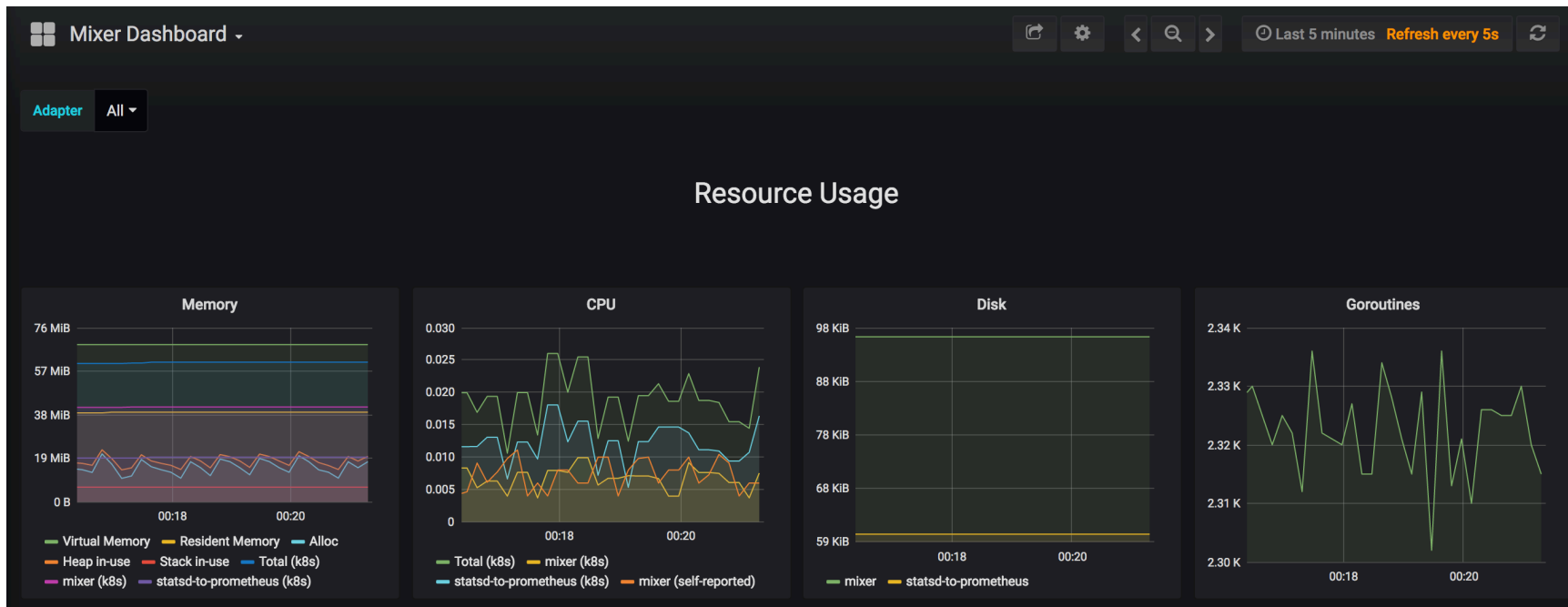


Telemetry



TLS certs to
Envoys

Telemetry





Security

<https://istio.io/docs/>

<https://www.ibm.com/cloud/info/istio>

Clusters / mycluster



mycluster

Expires in a month

● Requested

Terminal (Beta)

Kubernetes Dashboard



Access

Overview

Worker Nodes

Worker Pools

Add-ons



Managed Istio

Beta

Istio is a service mesh for providing a uniform way to integrate microservices, manage traffic flow across microservices, enforce policies and aggregate telemetry data. The Istio add-on simplifies the installation and maintenance of your istio control plane so you can focus on managing your microservices. [Learn More.](#)



Install

Labs

<http://ibm.biz/SSTIstioLab>

IBM