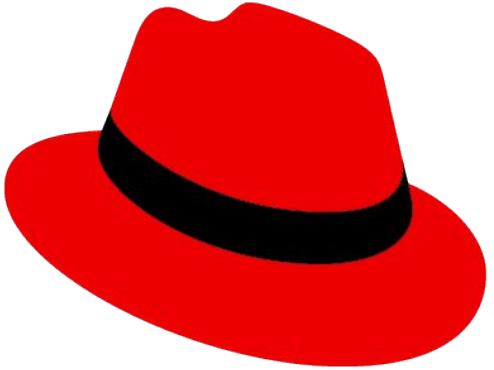


Kubernetes Patterns

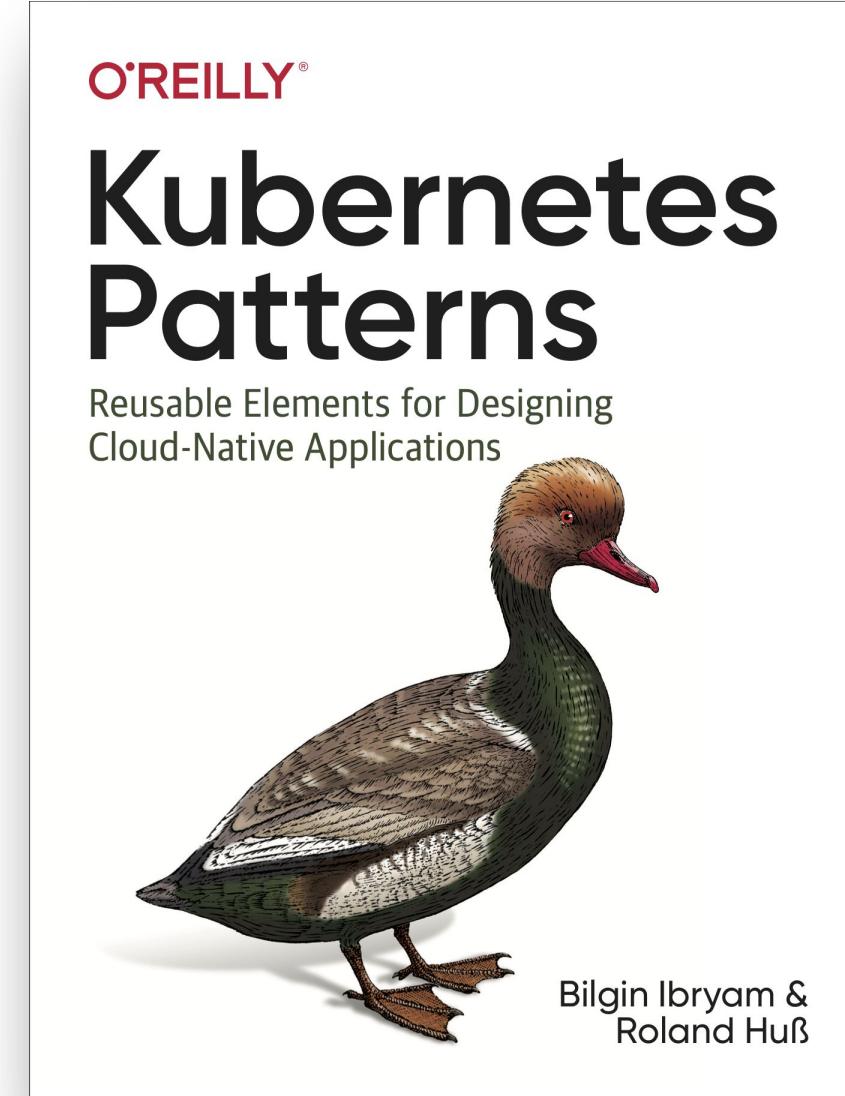
Reusable Elements for Designing
Cloud-Native Applications

Dr. Roland Huß
Principal Software Engineer
@ro14nd
<https://k8spatterns.io>



Red Hat

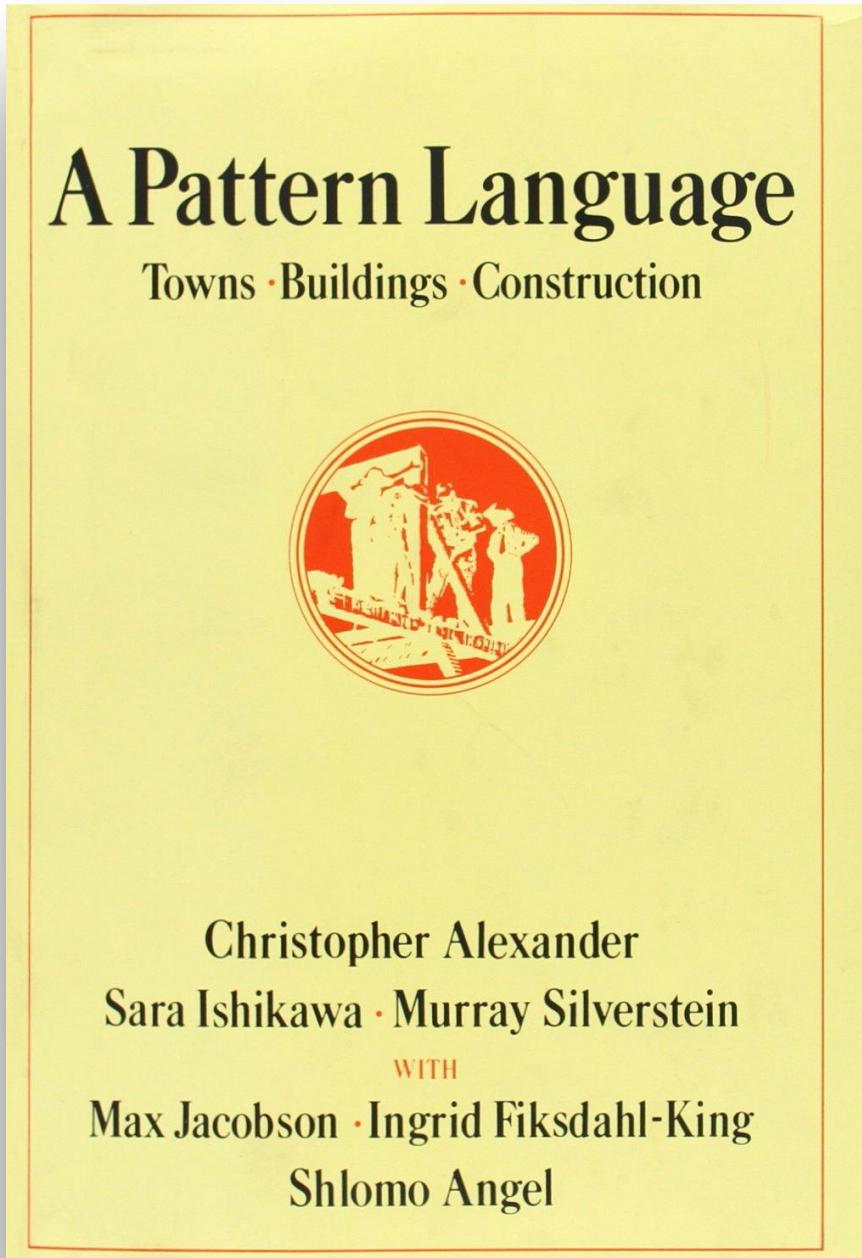
New Logo !



<https://k8spatterns.io>

The background of the image is a intricate geometric pattern. It features a repeating motif of overlapping diamond shapes, creating a sense of depth and movement. The colors used are earthy tones, including various shades of brown, tan, and black, which are rendered with a slightly textured, possibly marbled or wood-grain-like appearance.

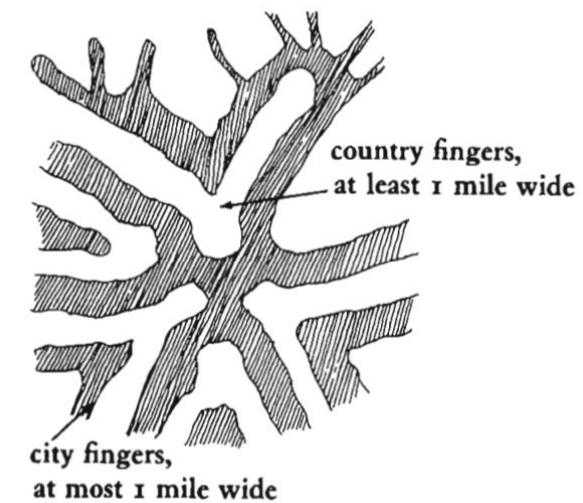
Patterns



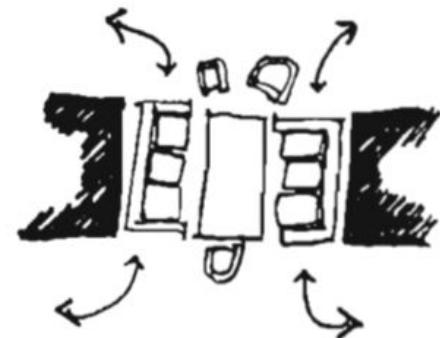
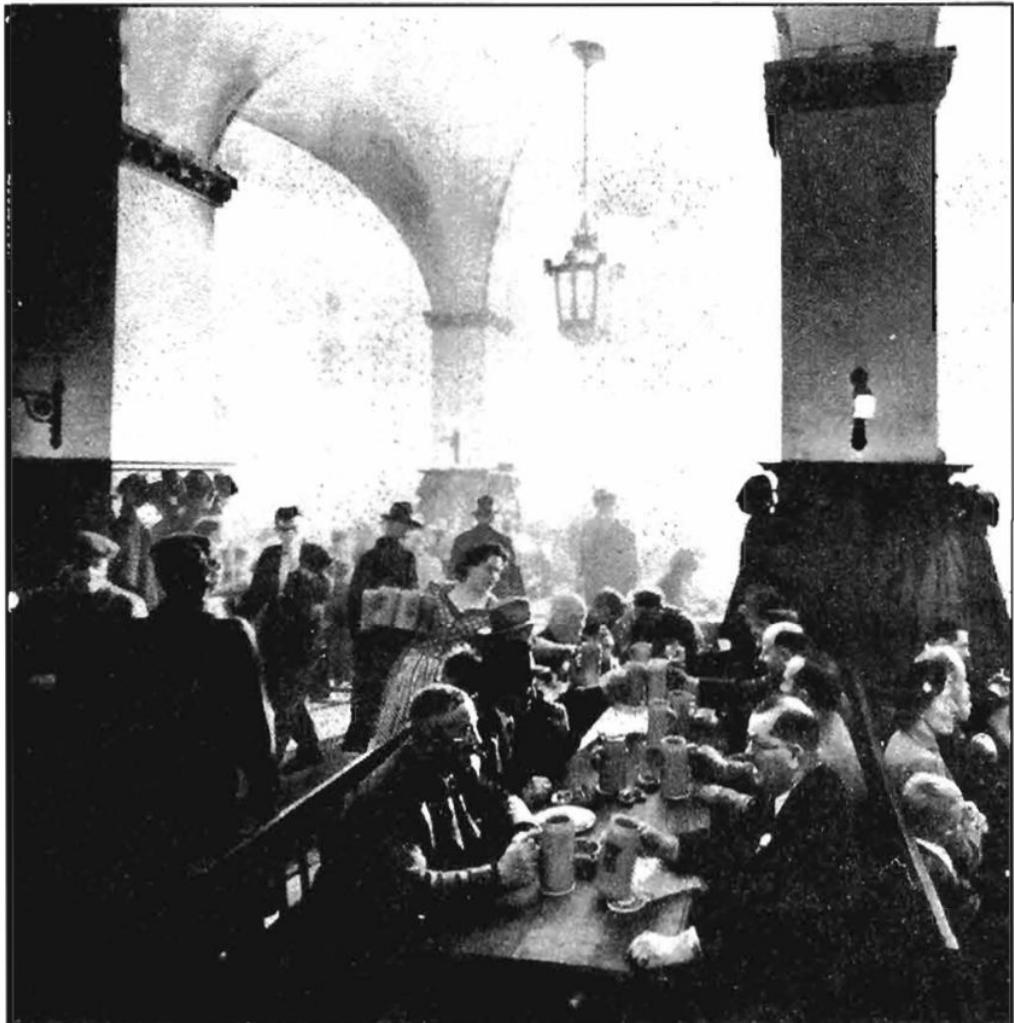
3 CITY COUNTRY FINGERS**



*When the countryside is far away
the city becomes a prison.*



90 BEER HALL



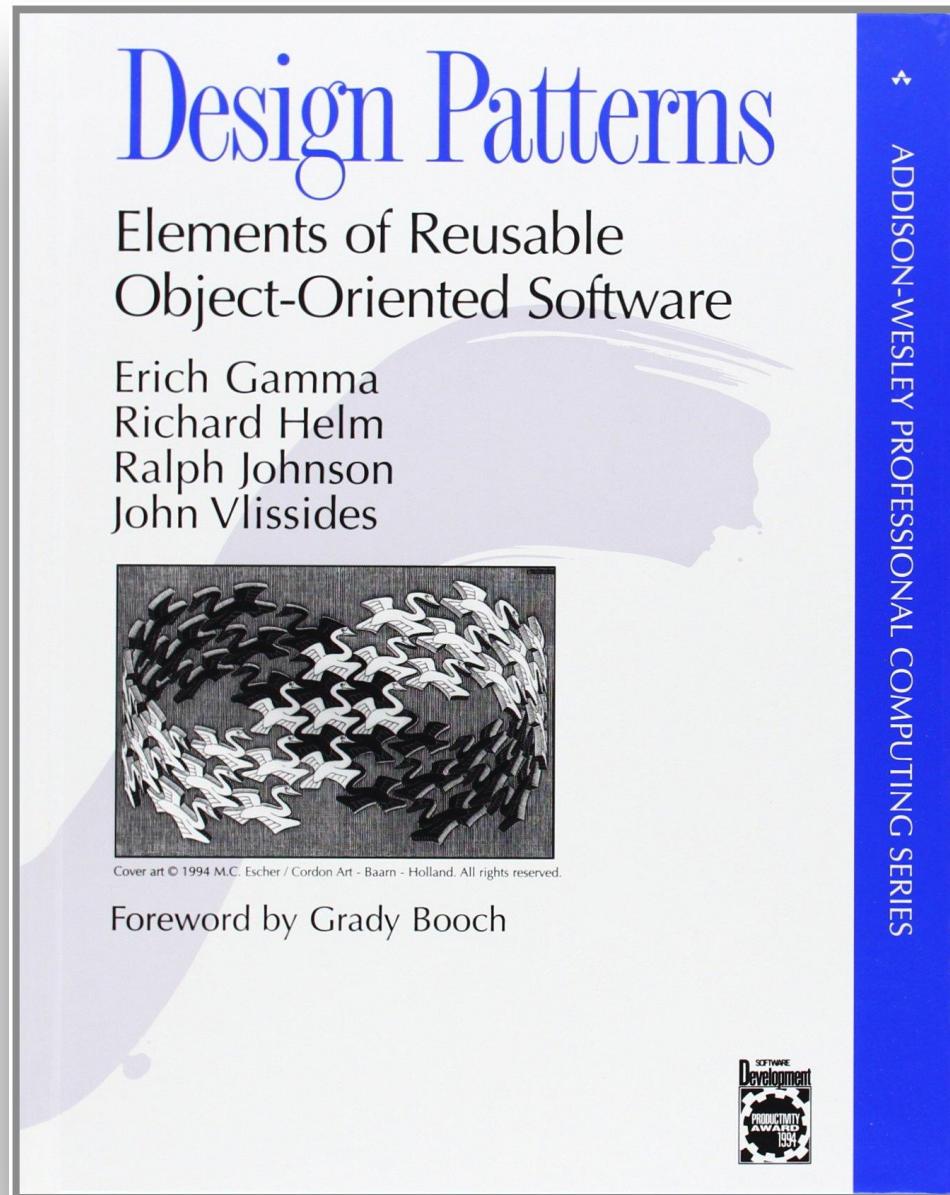
The open alcove—supports the fluidity of the scene.

criss-cross paths



activities

open alcoves



Patterns

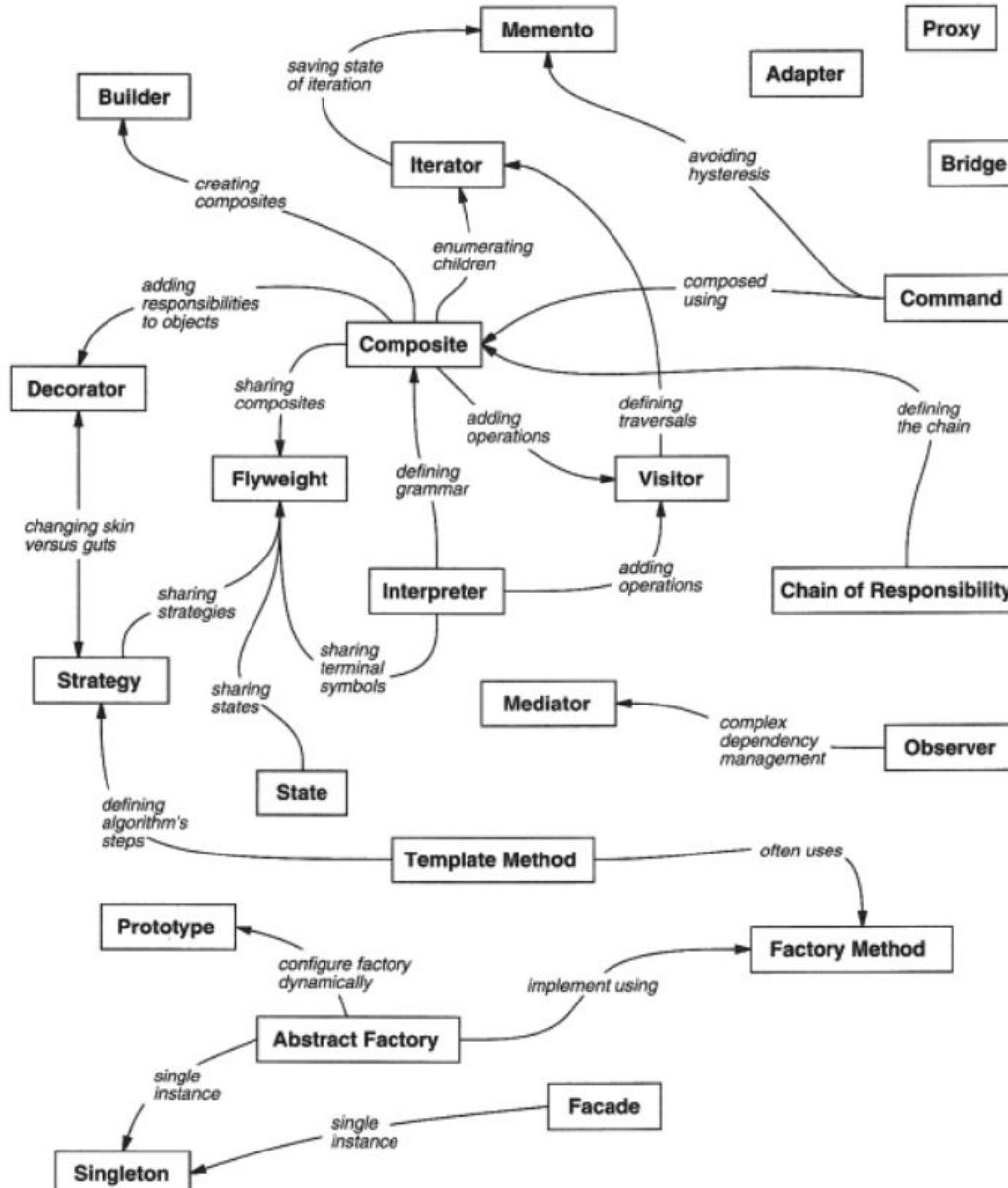
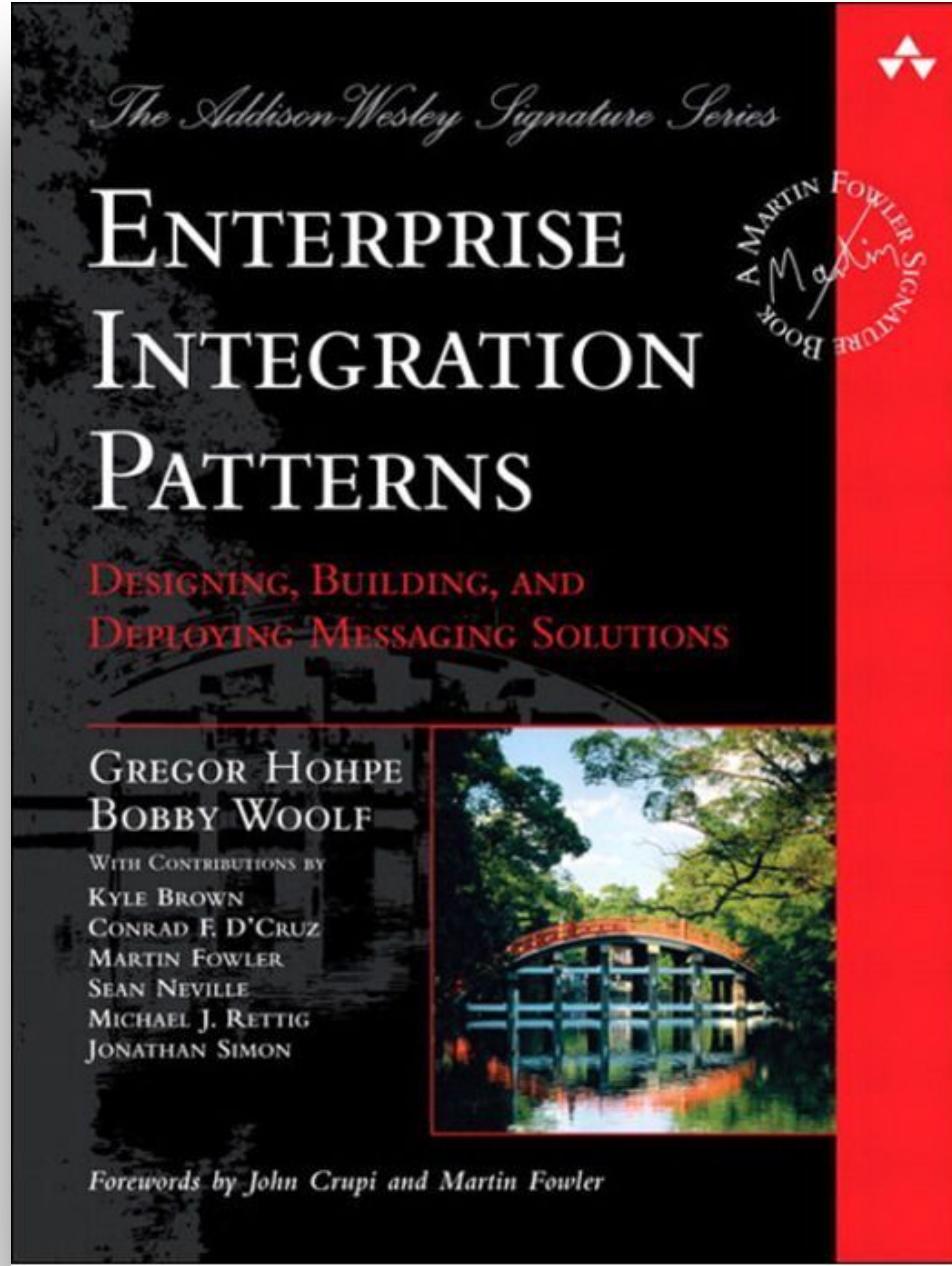


Figure 1.1: Design pattern relationships



A large, dark wooden ship's wheel is positioned in the foreground, angled towards the left. A rectangular metal plaque is attached to the right side of the wheel, featuring the words "U.S. COAST GUARD" in both English and Spanish. The background is a blurred photograph of a sunset or sunrise over the ocean, with warm orange and yellow hues reflecting on the water.

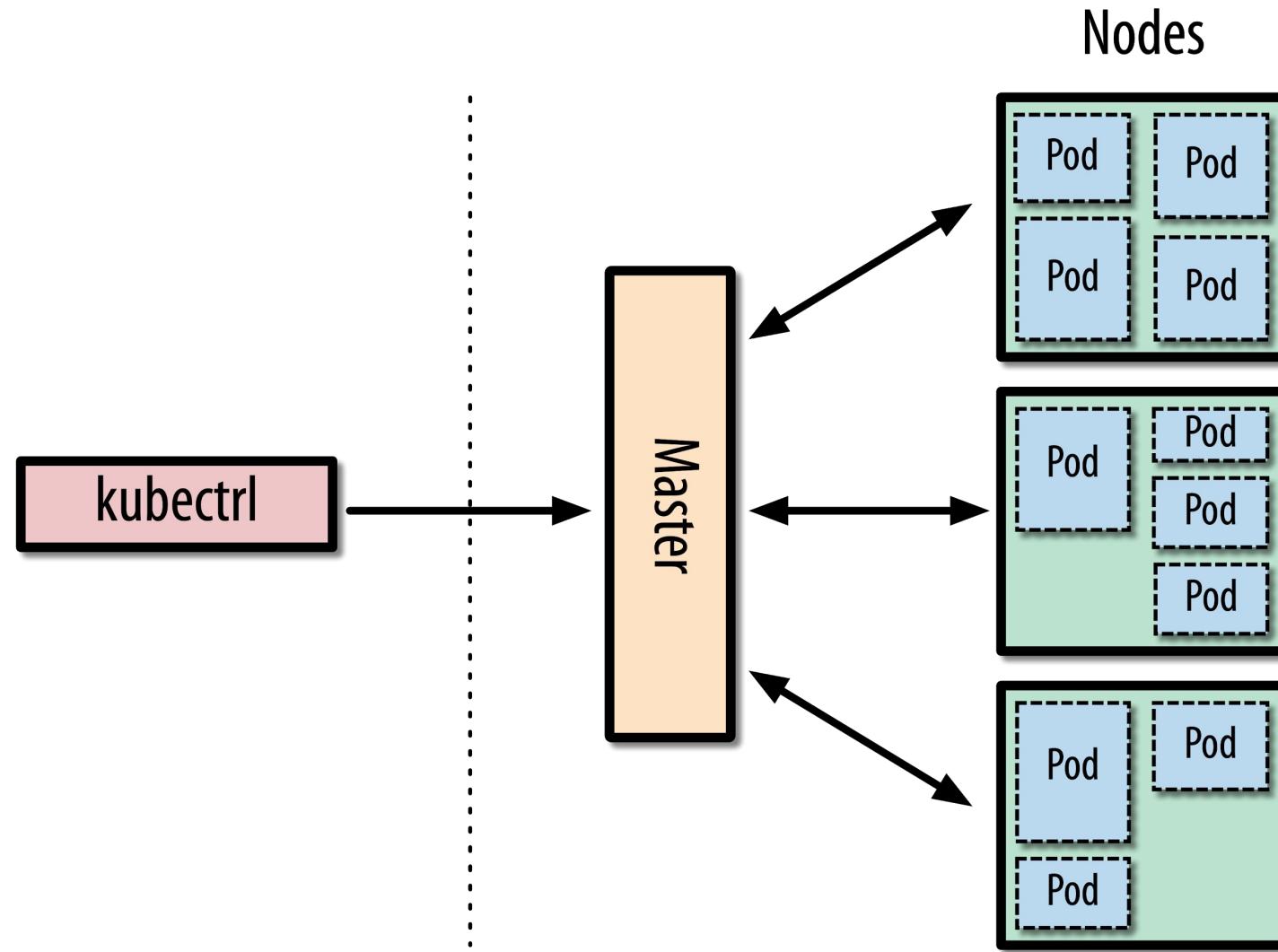
Kubernetes

Kubernetes

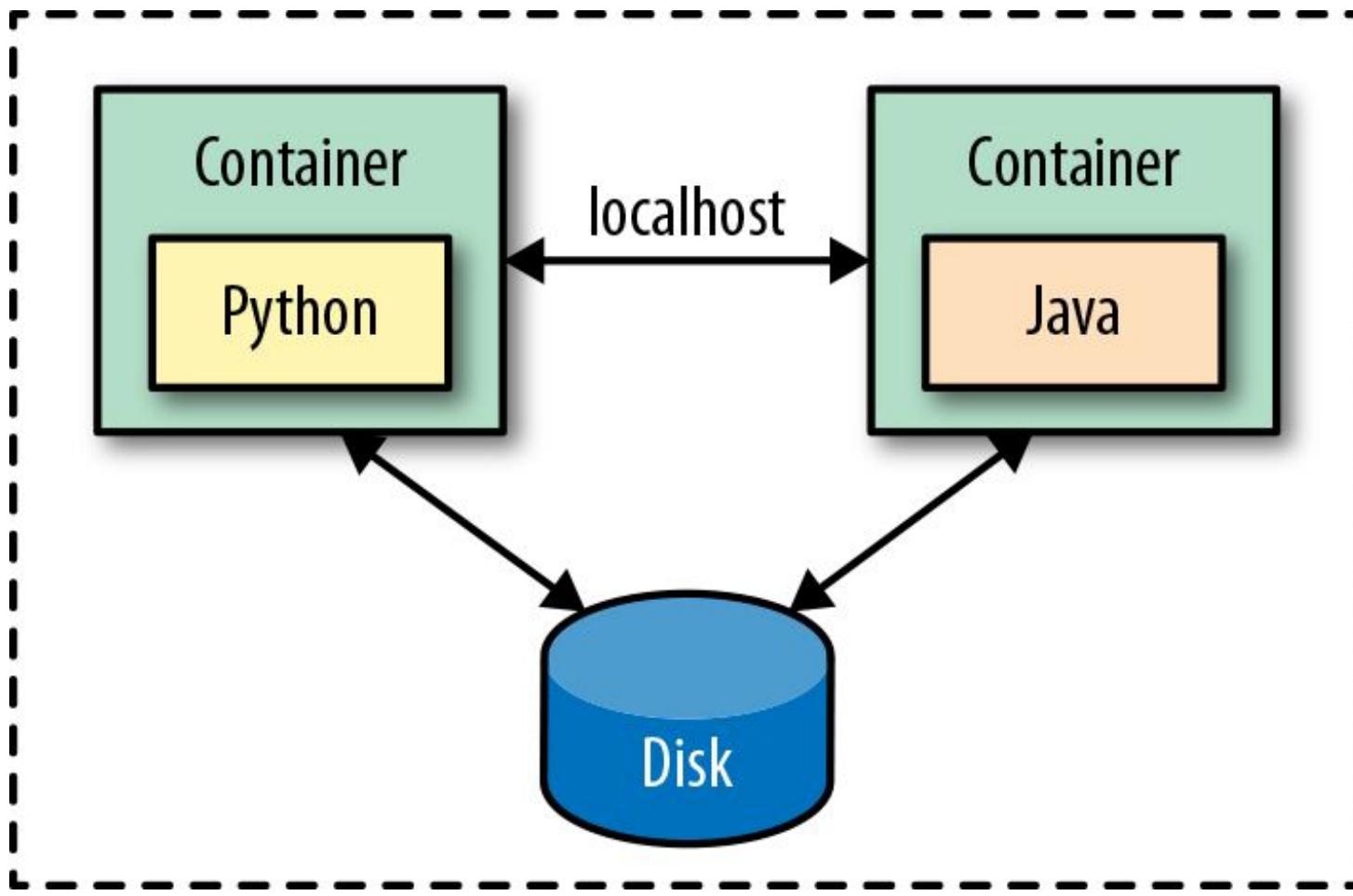


- Open Source container orchestration system started by Google in 2014
 - ※ Scheduling
 - ※ Self-healing
 - ※ Horizontal and vertical scaling
 - ※ Service discovery
 - ※ Rollout and Rollbacks
- Declarative resource-centric REST API

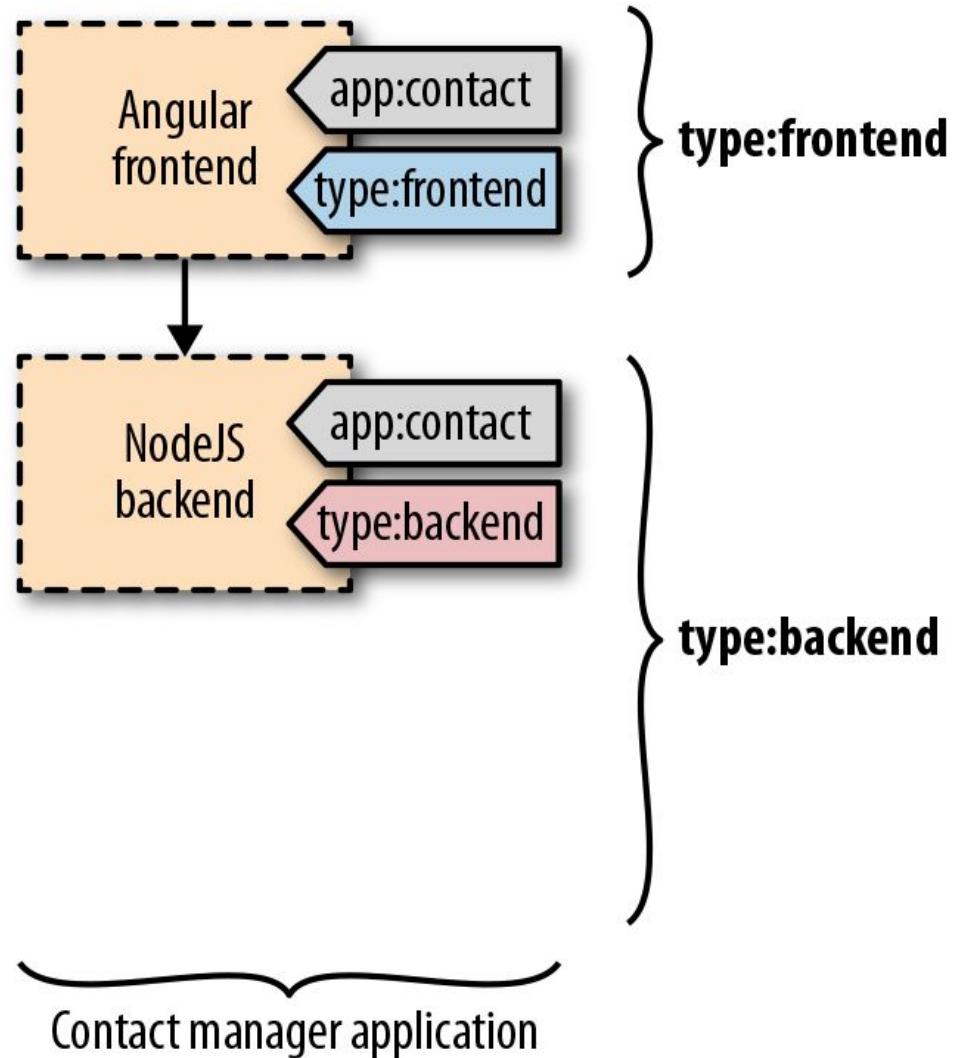
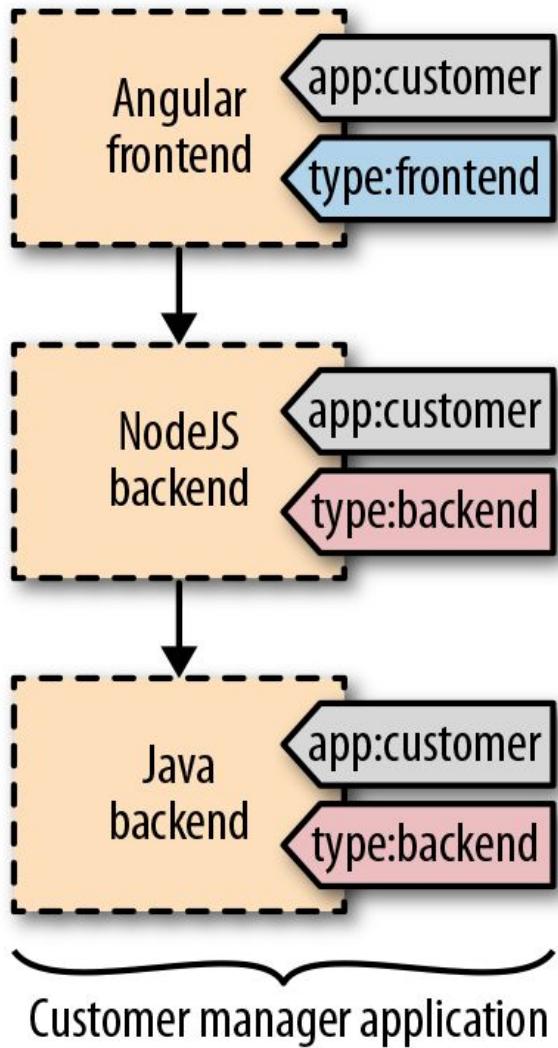
Architecture

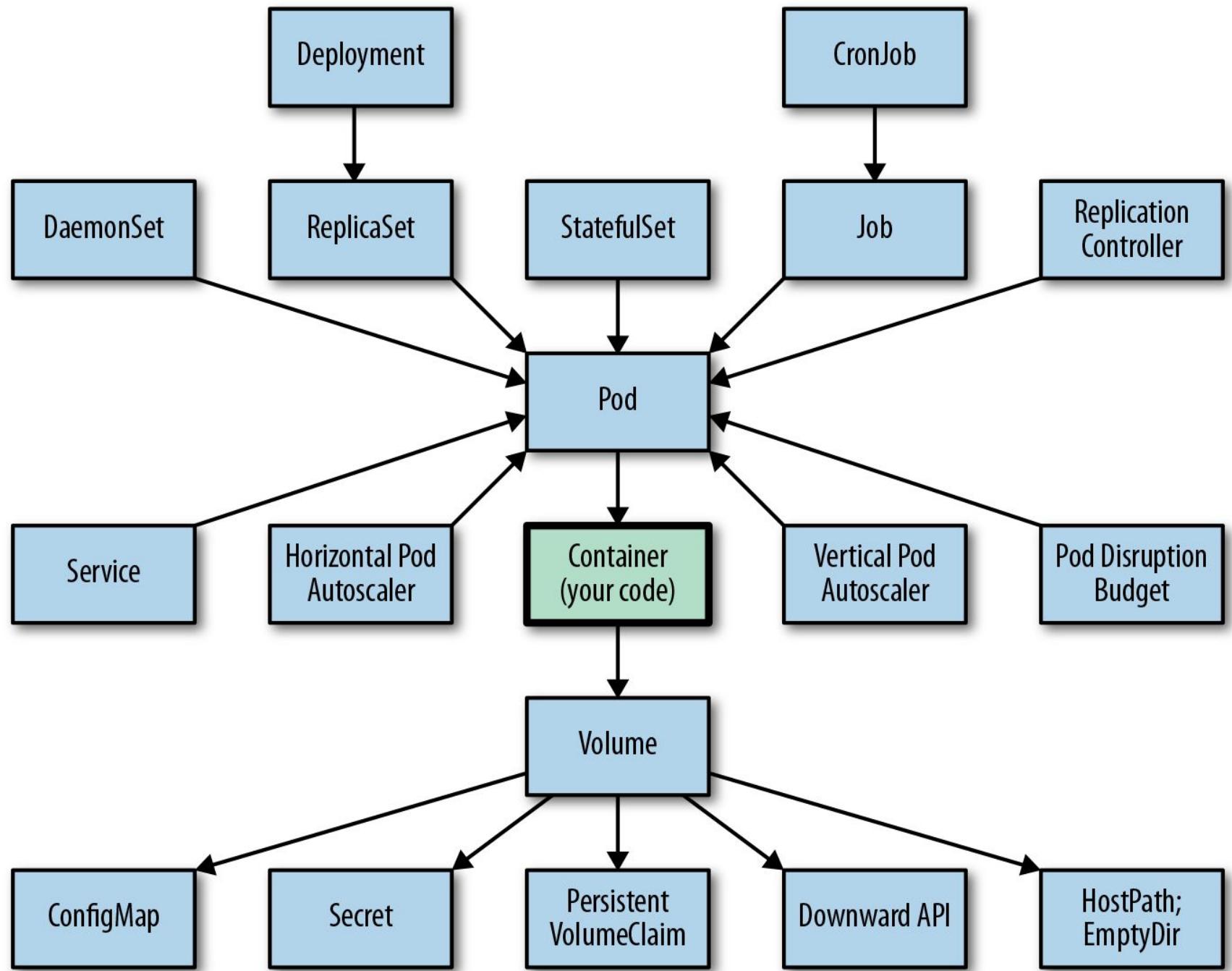


Pod



Labels







Predictable Demands

How to declare application requirements.

Application Requirements

- Declared requirements help in
 - Scheduling decisions
 - Capacity planning
 - Matching infrastructure services
- Runtime dependencies
 - Persistent Volumes
 - Host ports
 - Dependencies on ConfigMaps and Secrets

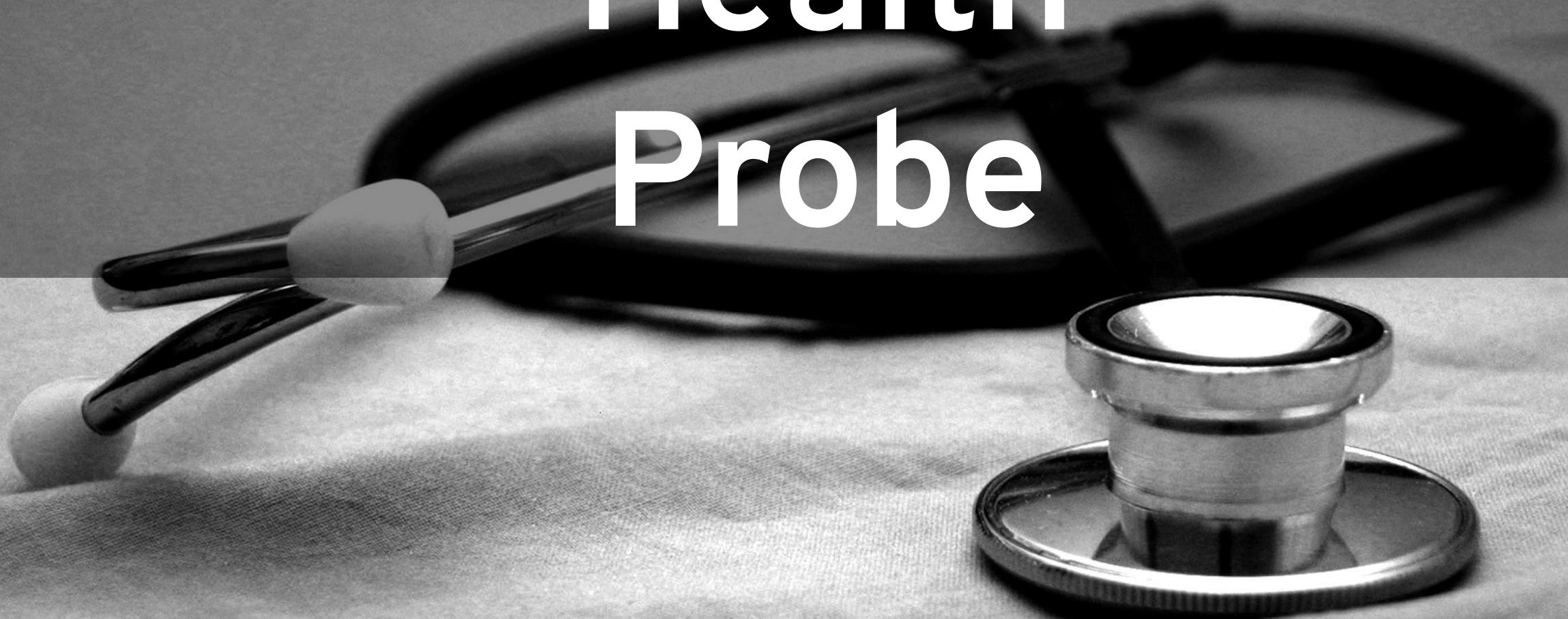
Resource Profile

```
apiVersion: v1
kind: Pod
metadata:
  name: http-server
spec:
  containers:
  - image: nginx
    name: nginx
  resources:
    requests:
      cpu: 200m
      memory: 100Mi
    limits:
      cpu: 300m
      memory: 200Mi
```

Quality-of-Service Classes

- **Best Effort**
 - No requests or limits
- **Burstable**
 - requests < limits
- **Guaranteed**
 - requests == limits

Health Probe



How to communicate an application's health state to Kubernetes.

Monitoring Health

- Process health checks
 - Checks running process
 - Restarts container if container process died
- Application provided Health Probes
 - **Liveness Probe:** Check application health
 - **Readiness Probe:** Check readiness to process requests

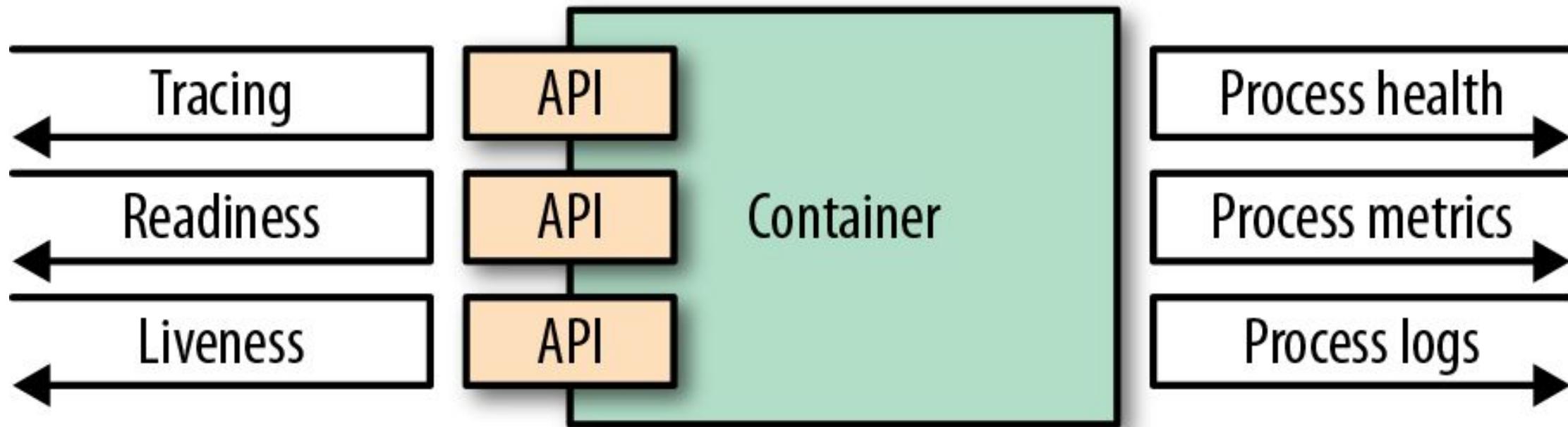
Liveness & Readiness

- Liveness Probe
 - Restarting containers if liveness probes fail
- Readiness Probe
 - Removing from service endpoint if readiness probe fails
- Probe methods
 - HTTP endpoint
 - TCP socket endpoint
 - Unix command's return value

Health Probe

```
apiVersion: v1
kind: Pod
metadata:
  name: pod-with-readiness-check
spec:
  containers:
  - image: k8spatterns/random-generator:1.0
    name: random-generator
    livenessProbe:
      httpGet:
        path: /actuator/health
        port: 8080
      initialDelaySeconds: 30
    readinessProbe:
      exec:
        command: [ "stat", "/var/run/random-generator-ready" ]
```

Container Observability Options

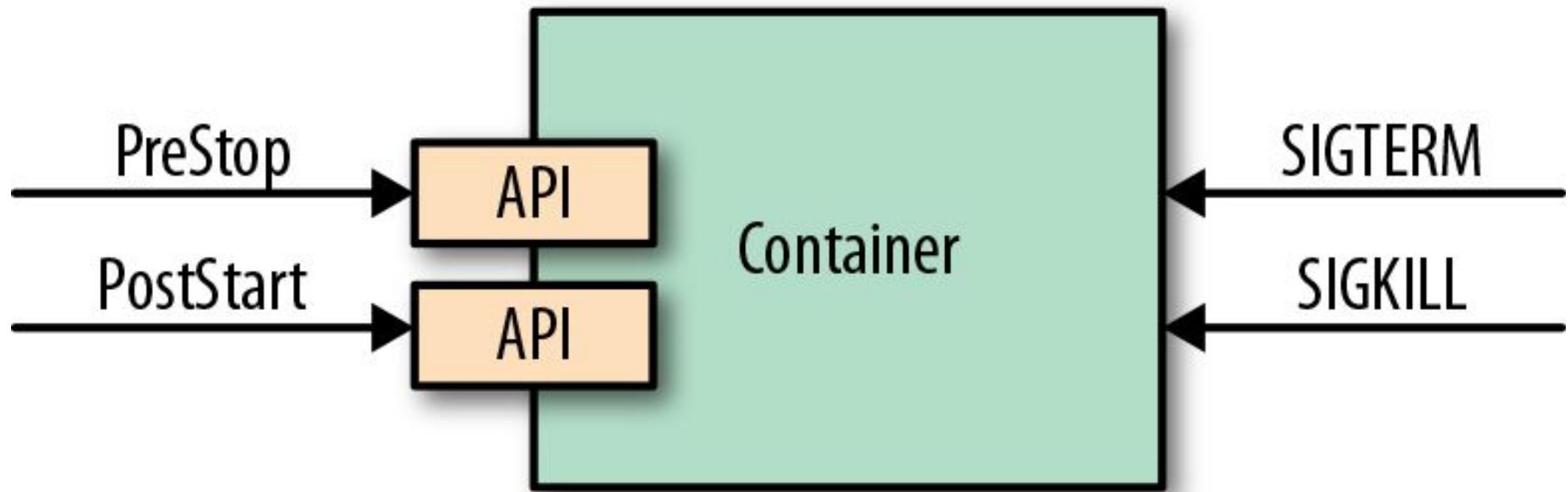




Managed Lifecycle

How applications should react on lifecycle events.

Managed Container Lifecycle



Lifecycle Events

- SIGTERM
 - Initial event issued when a container is going to shutdown
 - Application should listen to this event to cleanup properly and then exit
- SIGKILL
 - Final signal sent after a grace period which can't be catched
 - terminationGracePeriodSeconds: Time to wait after SIGTERM (default: 30s)

Lifecycle Hooks

- **postStart**
 - Called after container is created
 - Runs in parallel to the main container
 - Keeps Pod in status *Pending* until exited successfully
 - **exec** or **httpGet** handler types (like *Health Probe*)
- **preStop**
 - Called before container is stopped
 - Same purpose & semantics as for SIGTERM

Managed Lifecycle

```
apiVersion: v1
kind: Pod
metadata:
  name: pre-stop-hook
spec:
  containers:
  - image: k8spatterns/random-generator:1.0
    name: random-generator
    lifecycle:
      postStart:
        exec:
          command:
          - sh
          - -c
          - sleep 30 && echo "Wake up!" > /tmp/postStart_done
      preStop:
        httpGet:
          port: 8080
          path: /shutdown
```

The background of the image consists of a dense, repeating pattern of small, silver-colored cylindrical containers, likely aerosol cans, arranged in a grid. They are positioned on a dark blue surface. The lighting creates highlights on the metallic surfaces of the containers.

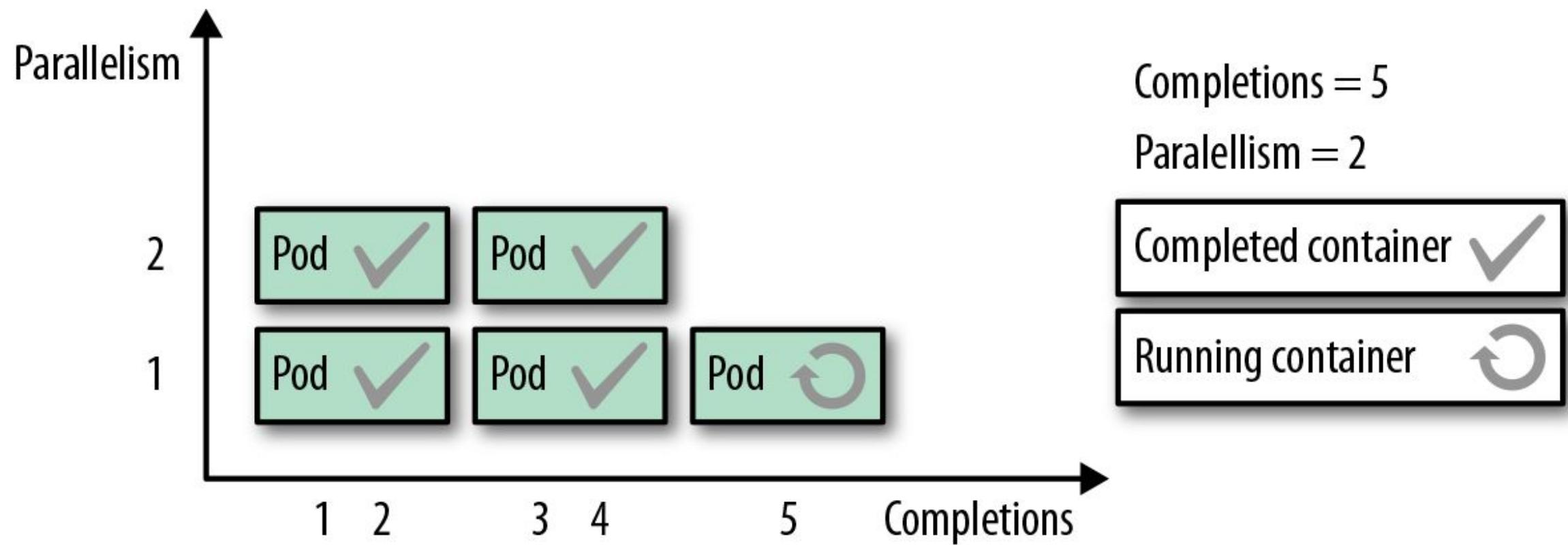
Batch Job

How to run short-lived Pods reliably until completion.

Job

- Resource for a predefined *finite* unit-of-work
- Survives cluster restarts and node failures
- Support for multiple Pod runs
- `.spec.completions`
 - How many Pods to run to complete Job
- `.spec.parallelism`
 - How many Pods to run in parallel

Parallel Batch Job



Job Types

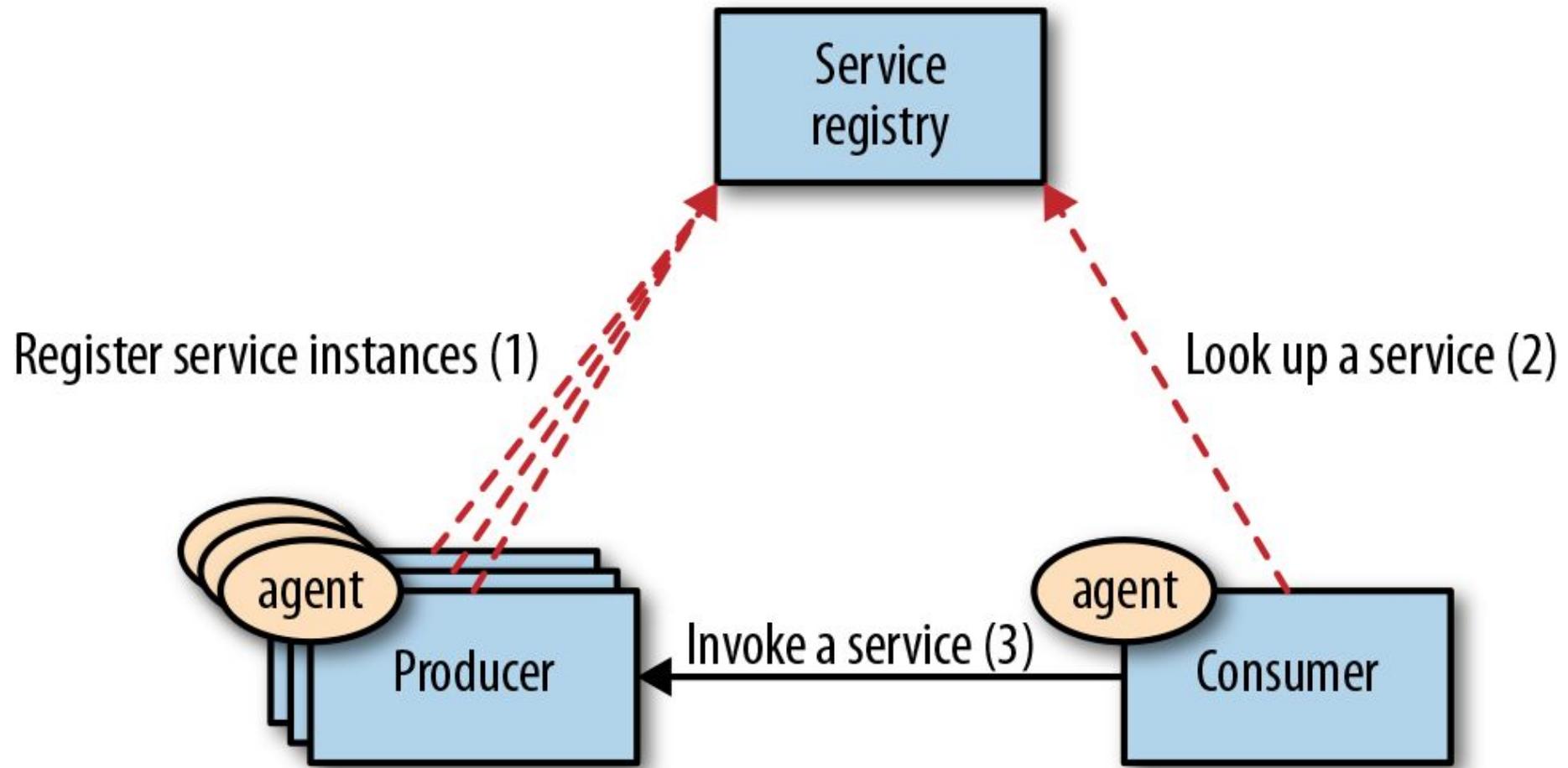
- Single Pod Job
 - completions = 1 and parallelism = 1
- Fixed completion count Jobs
 - completions > 1
 - One Pod per work item
- Work queue Jobs
 - completions = 1 and parallelism > 1
 - All Pods need to finish, with at least one Pod exiting successfully
 - Pods needs to coordinate to shutdown in a coordinate fashion



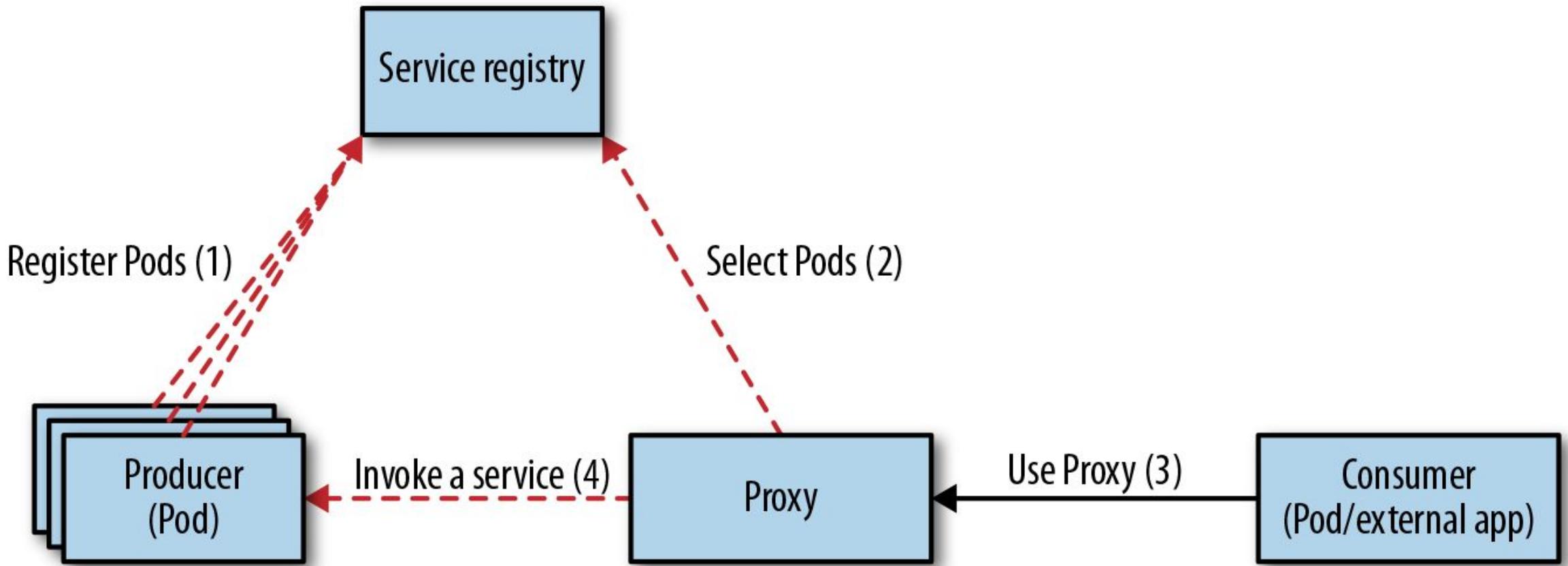
Service Discovery

How to discover and use services.

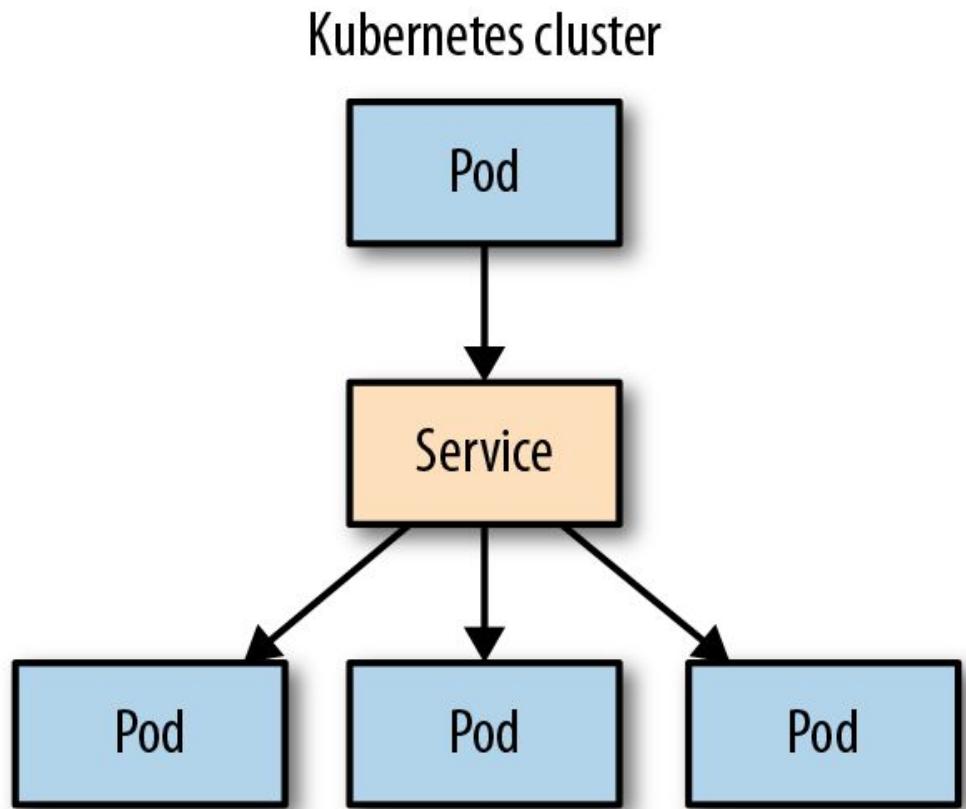
Client-side Service Discovery (non Kubernetes)



Server-side Service Discovery (Kubernetes)

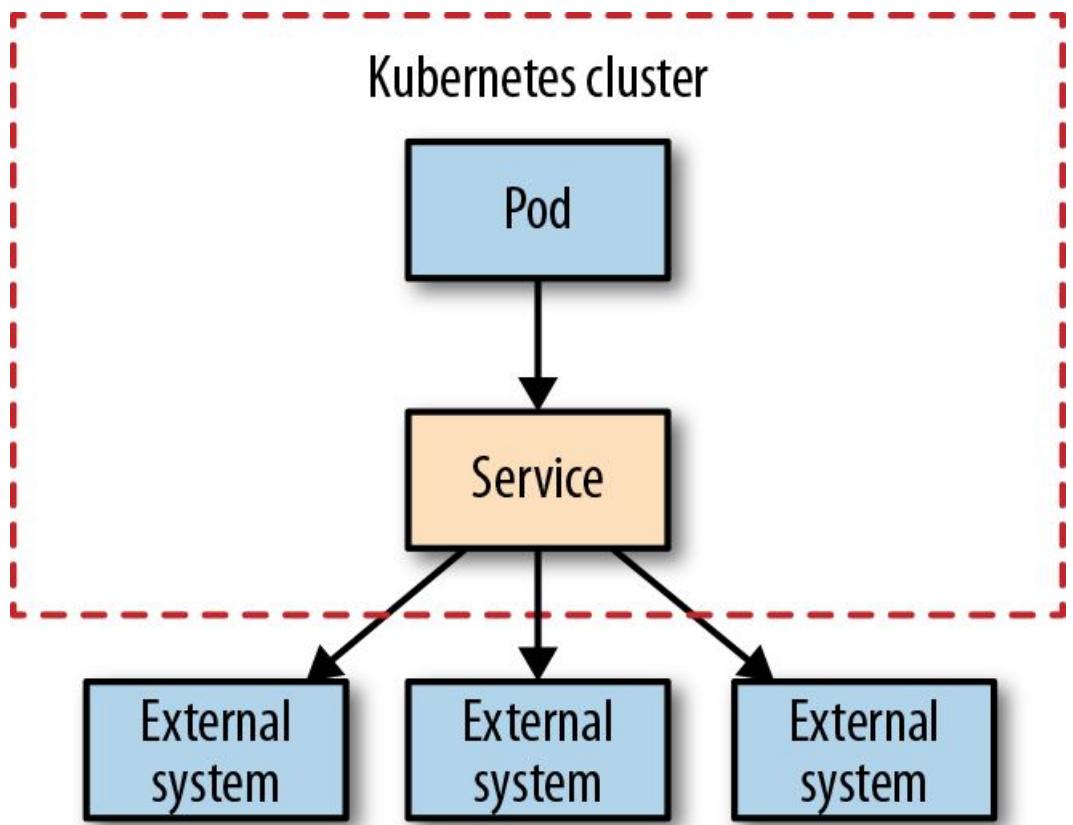


Internal Service Discovery



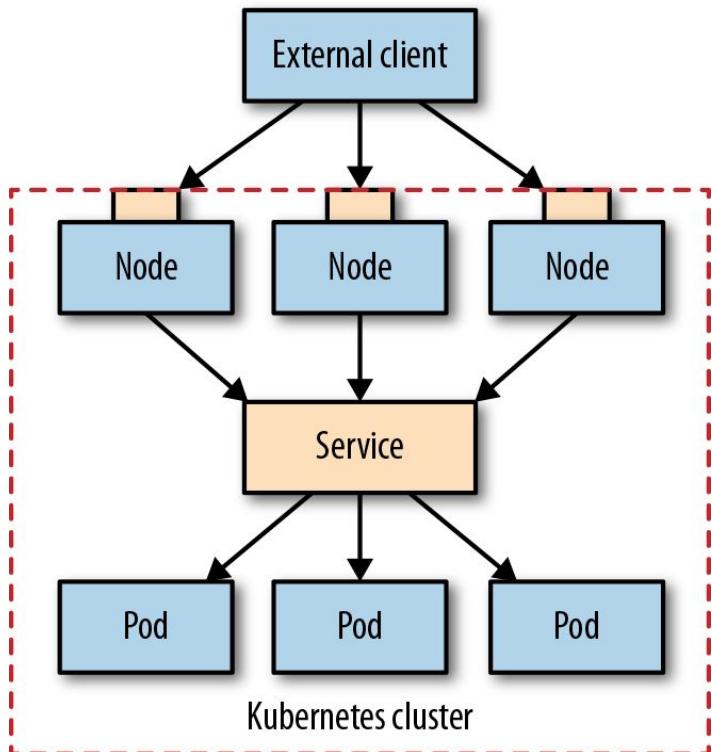
- Discovery through DNS lookups
- Pods picked by label selector
- Multiple ports per Service
- Session affinity on IP address
- Successful readiness probes required for routing
- Virtual IP address for each Service

Manual Service Discovery

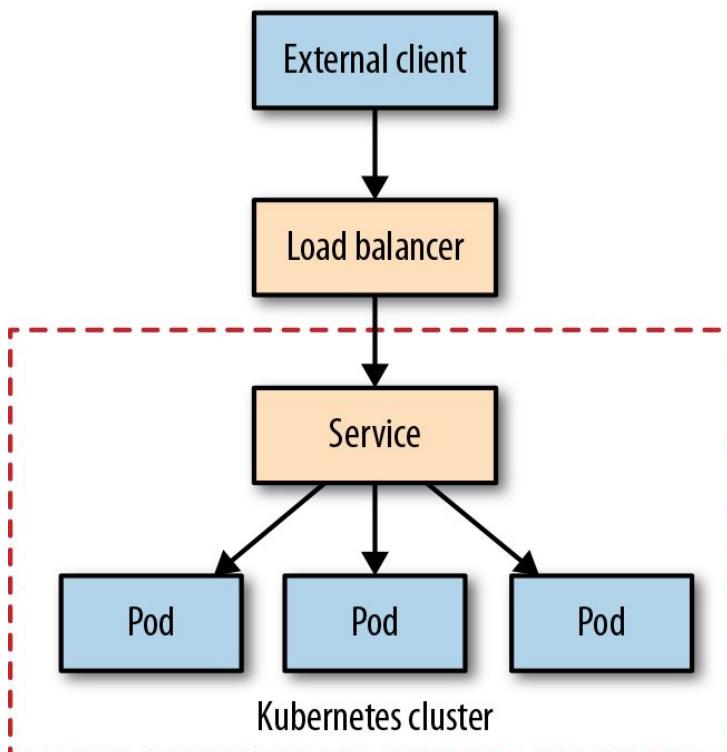


- Service without selector
- Manually creating Endpoint resource with the same name as the Service
- Service of type **ExternalName** map are registered as DNS CNAMEs

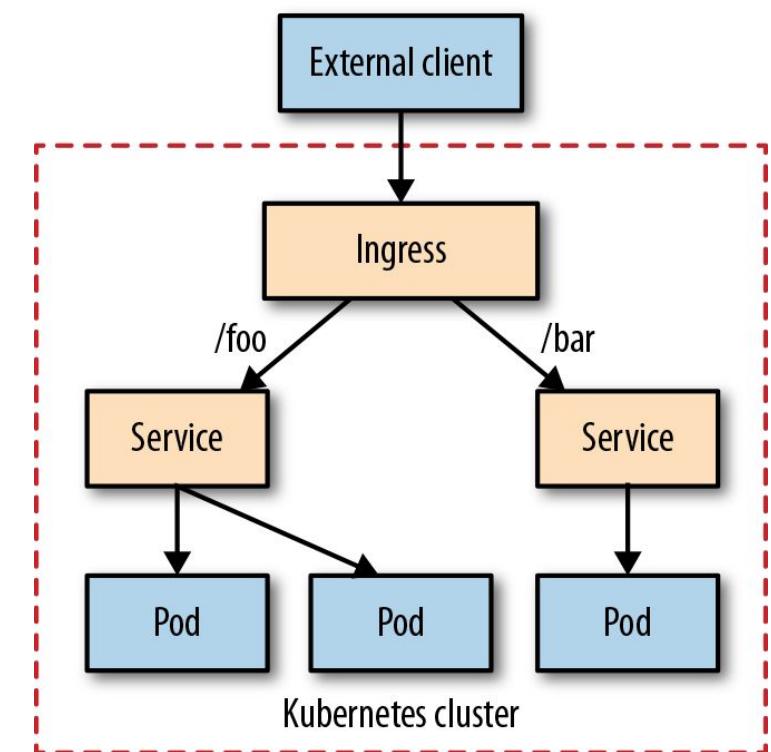
Node Port



Load Balancer



Ingress

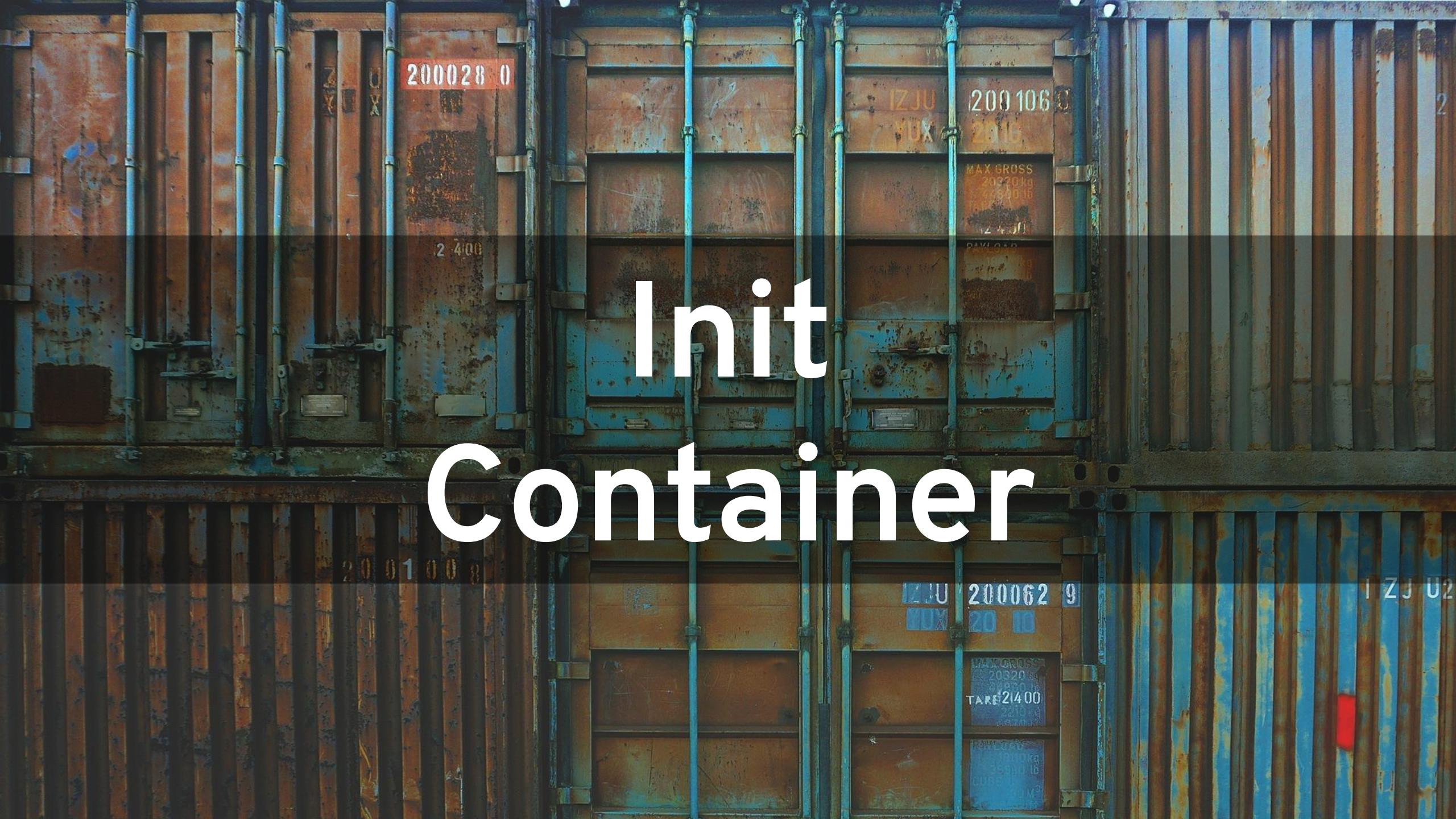


Ingress

```
apiVersion: extensions/v1beta1
kind: Ingress
metadata:
  name: random-generator
spec:
  rules:
  - http:
    paths:
    - path: /
      backend:
        serviceName: random-generator
        servicePort: 8080
    - path: /cluster-status
      backend:
        serviceName: cluster-status
        servicePort: 80
```

Service Discovery

Name	Configuration	Client type	Summary
ClusterIP	<code>type: ClusterIP</code> <code>.spec.selector</code>	Internal	The most common internal discovery mechanism
Manual IP	<code>type: ClusterIP</code> <code>kind: Endpoints</code>	Internal	External IP discovery
Manual FQDN	<code>type: ExternalName</code> <code>.spec.externalName</code>	Internal	External FQDN discovery
Headless Service	<code>type: ClusterIP</code> <code>.spec.clusterIP: None</code>	Internal	DNS-based discovery without a virtual IP
NodePort	<code>type: NodePort</code>	External	Preferred for non-HTTP traffic
LoadBalancer	<code>type: LoadBalancer</code>	External	Requires supporting cloud infrastructure
Ingress	<code>kind: Ingress</code>	External	L7/HTTP-based smart routing mechanism



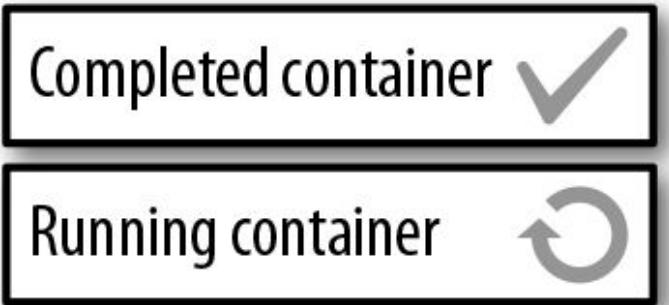
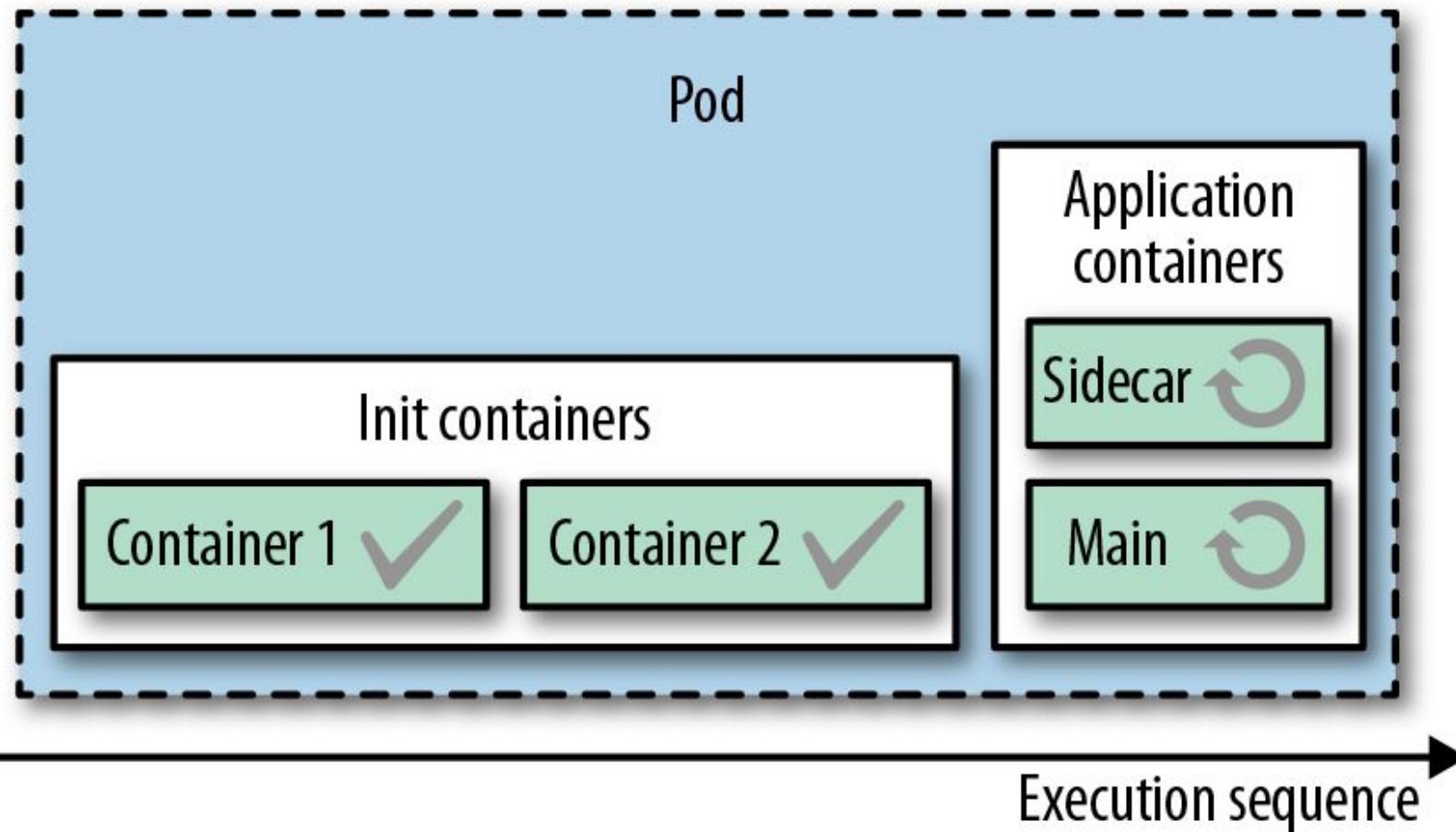
Init Container

How to separate
initialization steps from
the main application.

Init Container

- Part of a Pod
- One shot actions before application starts
- Needs to be idempotent
- Has own resource requirements

Init Container



Init Container

```
apiVersion: v1
kind: Pod
.....
spec:
  initContainers:
  - name: download
    image: axeclbr/git
    command: [ "git", "clone", "https://github.com/myrepo", "/data" ]
    volumeMounts:
      - mountPath: /var/lib/data
        name: source
  containers:
  - name: run
    image: docker.io/centos/httpd
    volumeMounts:
      - mountPath: /var/www/html
        name: source
  volumes:
  - emptyDir: {}
    name: source
```

The image is a collage of three photographs related to sidecar racing. The top photograph shows a rider in a purple and white patterned suit and helmet riding a sidecar through mud, with another rider visible behind them. The middle photograph is a close-up of a sidecar, heavily covered in mud, with its number plate partially visible. The bottom photograph shows a rider in a dark suit and helmet leaning into a turn on a muddy track.

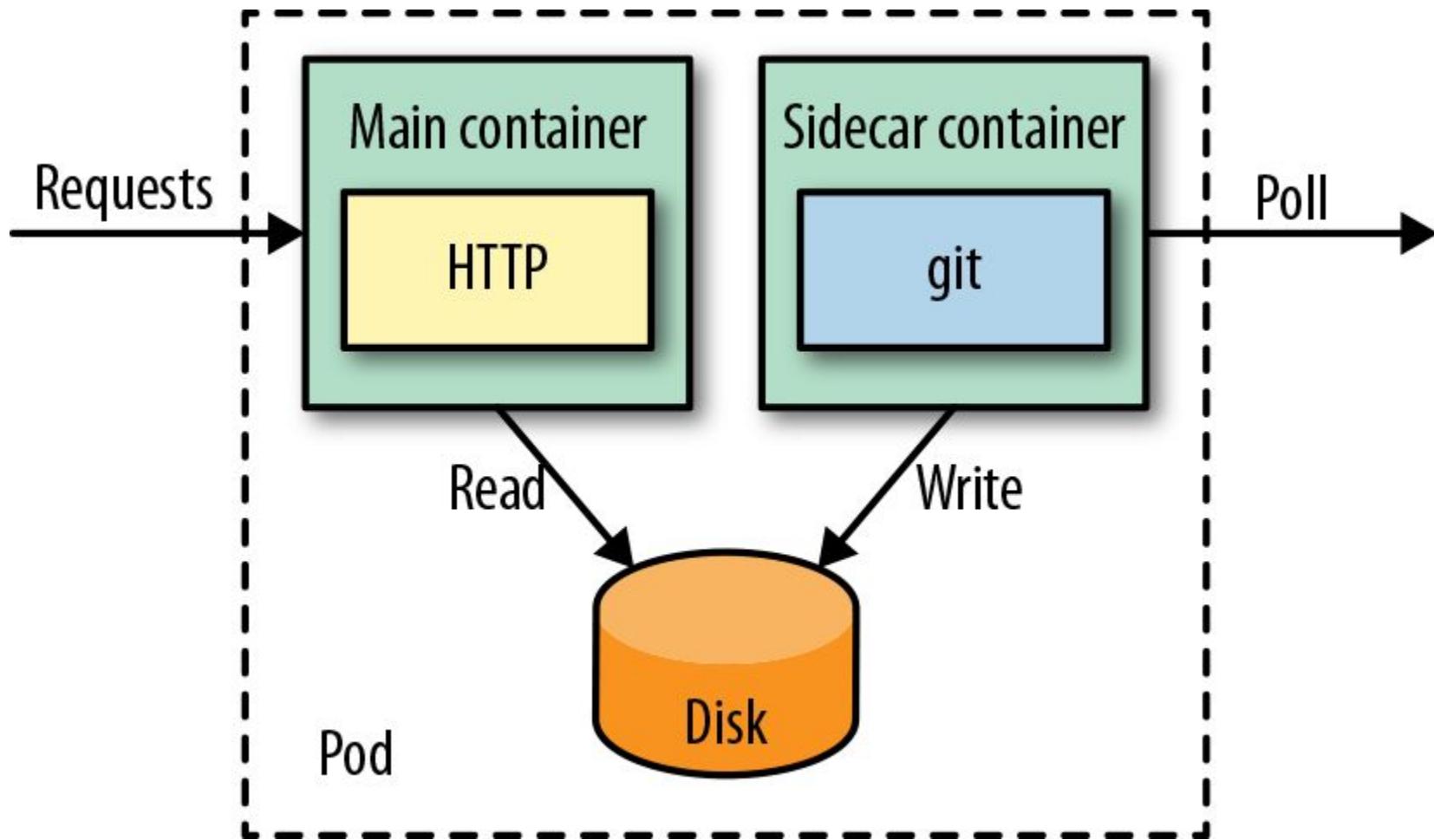
Sidecar

How to enhance the
functionality of an application
without changing it.

Sidecar

- Runtime collaboration of containers
- Connected via shared resources:
 - Network
 - Volumes
- Similar what AOP is for programming
- Separation of concerns

Sidecar

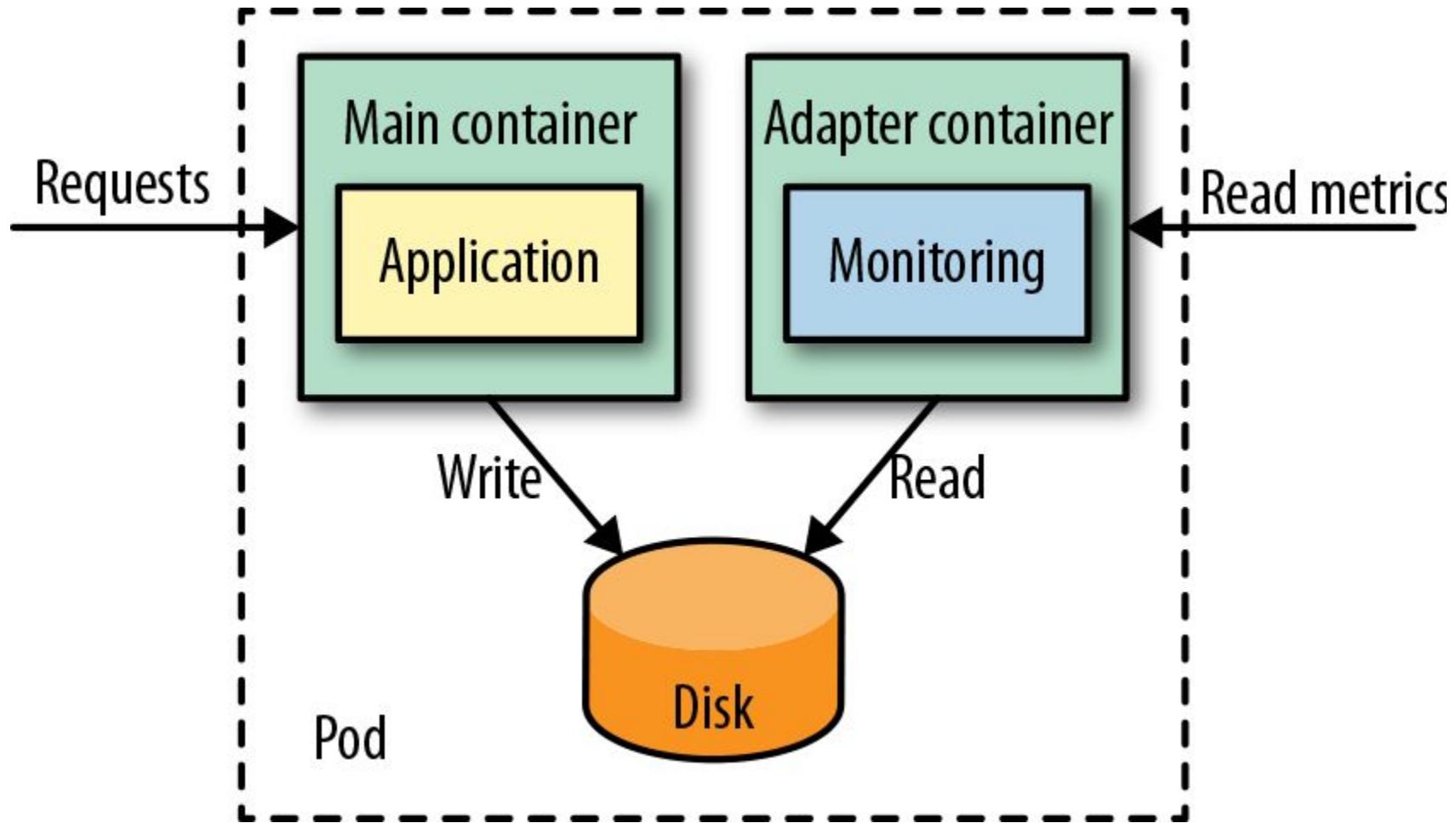


Adapter



How to decouple access to a
container **from** the outside
world.

Adapter

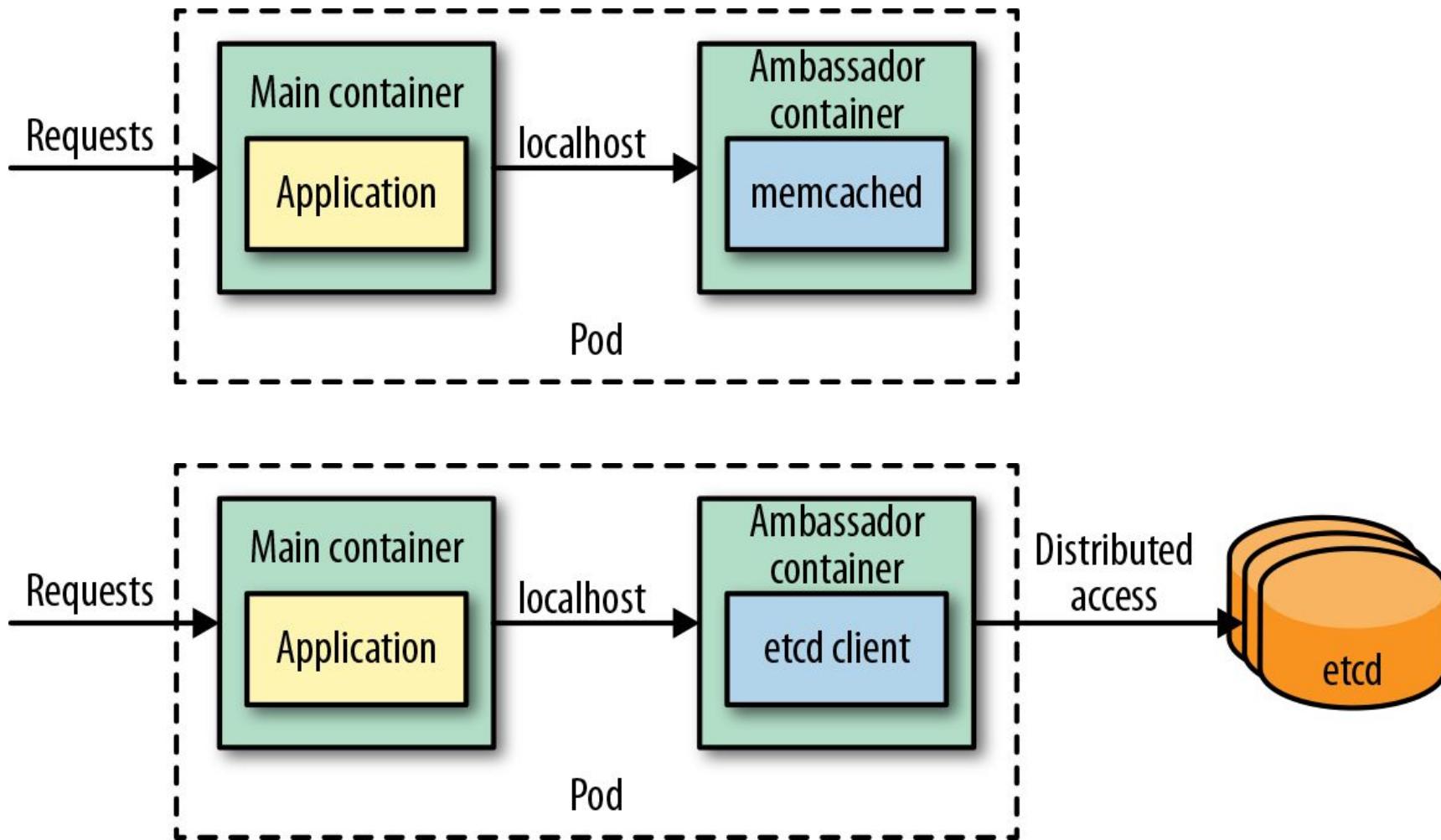




Ambassador

How to decouple a container's access **to** the outside world

Ambassador



Controller



How to get from the
current state to a declared
target state.

State Reconciliation

- Kubernetes as distributed state manager
- Make the **actual** state more like the declared **target** state.



Observe - Discover the actual state

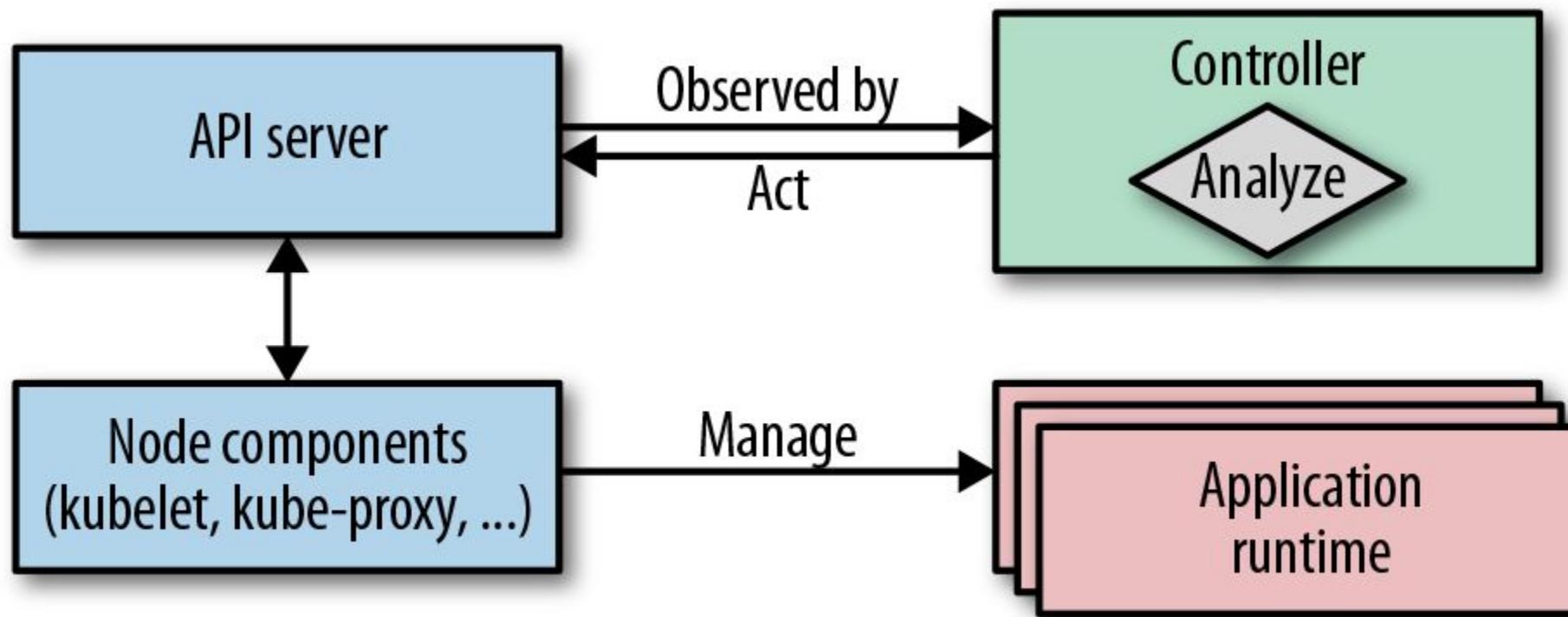


Analyze - Determine difference to target state



Act- Perform actions to drive the actual to the desired state

Observe - Analyze - Act



Common Triggers

- Labels
 - Indexed by backend
 - Suitable for selector-like functionality
 - Limitation on charset for names and values
- Annotations
 - No syntax restrictions
 - Not indexed
- ConfigMaps
 - Good for complex structured state declarations
 - Simple alternative to CustomResourceDefinitions

ConfigMap Watch Controller

```
namespace=${WATCH_NAMESPACE:-default}
base=http://localhost:8001
ns=namespaces/$namespace

curl -N -s $base/api/v1/${ns}/configmaps?watch=true | \
while read -r event
do
    type=$(echo "$event" | jq -r '.type')
    config_map=$(echo "$event" | jq -r '.object.metadata.name' )
    annotations=$(echo "$event" | jq -r '.object.metadata.annotations' )

    if [ $type = "MODIFIED" ]; then
        # Restart Pods using this ConfigMap
        # ...
    fi
done
```



Operator

How to encapsulate
operational knowledge
into executable software.

Definition

“” An **operator** is a Kubernetes **controller** that understands two domains: Kubernetes and *something else*. By combining knowledge of both areas, it can **automate** tasks that usually require a human operator that understands both domains.

Jimmy Zelinskie

<http://bit.ly/2Fjlx1h>

Technical:

Operator = Controller + CustomResourceDefinition

CustomResourceDefinition

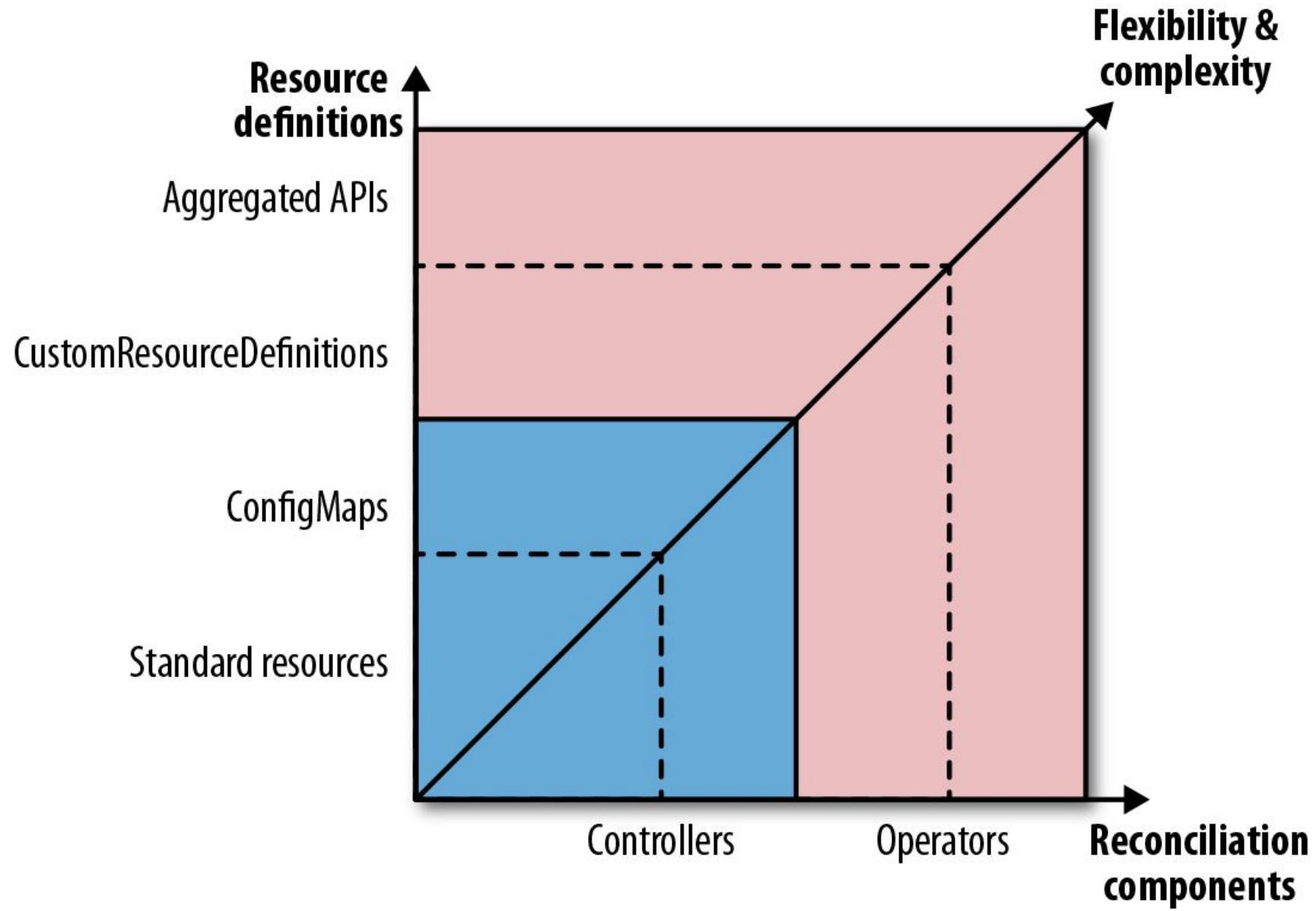
Custom resource is modelling a custom domain and managed through the Kubernetes API

```
apiVersion: apiextensions.k8s.io/v1beta1
kind: CustomResourceDefinition
metadata:
  name: configwatchers.k8spatterns.io
spec:
  scope: Namespaced
  group: k8spatterns.io
  version: v1
  names:
    kind: ConfigWatcher
    plural: configwatchers
  validation:
    openAPIV3Schema:
      ...

```

Custom Resource

```
kind: ConfigWatcher
apiVersion: k8spatterns.io/v1
metadata:
  name: webapp-config-watcher
spec:
  configMap: webapp-config
  podSelector:
    app: webapp
```



CRD Classification

- Installation CRDs
 - Installing and operating applications
 - Backup and Restore
 - Monitoring and self-healing
 - Example: Prometheus for installing Prometheus & components
- Application CRDs
 - Application specific domain concepts
 - Example: ServiceMonitor for registering Kubernetes service to be scraped by Prometheus

Operator Hub

The screenshot shows the OperatorHub.io website interface. At the top, there is a navigation bar with a search bar labeled "Search OperatorHub..." and a "Contribute" button. Below the header, a main banner reads "Welcome to OperatorHub.io" and "OperatorHub.io is a new home for the Kubernetes community to share Operators. Find an existing Operator or list your own today." On the left side, there are two filter sections: "PROVIDER" and "CAPABILITY LEVEL". The "PROVIDER" section has checkboxes for various providers, with "Jaeger" and "Red Hat" checked. The "CAPABILITY LEVEL" section has checkboxes for "Basic Install", "Seamless Upgrades", and "Full Lifecycle". In the center, there is a grid of operator cards. The first card is for "Jaeger Tracing", provided by Jaeger, which offers tracing, monitoring, and troubleshooting for microservices-based applications. The second card is for "Kubernetes Federation", provided by Red Hat, which enables hybrid cloud capabilities between clusters. The third card is for "MongoDB", provided by MongoDB, Inc., which facilitates easy deployments of MongoDB. The fourth card is for "Prometheus Operator", provided by Red Hat, which provides monitoring definitions for Kubernetes. The fifth card is for "Strimzi Kafka", provided by Red Hat, which runs an Apache Kafka cluster, including Kafka Connect, ZooKeeper, and more.

Provider

- Amazon Web Services (1)
- CNCF (1)
- Couchbase (1)
- Crunchy Data Solutions (1)
- Dynatrace (1)
- Jaeger (1)
- MongoDB (1)
- Percona (1)
- PlanetScale (1)
- Red Hat (3)
- Redis Labs (1)

Show less

Capability Level

- Basic Install (4)
- Seamless Upgrades (3)
- Full Lifecycle (6)

5 ITEMS

VIEW grid ▾ SORT A-Z ▾

 Jaeger Tracing provided by Jaeger Provides tracing, monitoring and troubleshooting microservices-based	 Kubernetes Federation provided by Red Hat Gain Hybrid Cloud capabilities between your clusters with Kubernetes Federation.	 MongoDB provided by MongoDB, Inc. The MongoDB Enterprise Kubernetes Operator enables easy deploys of MongoDB.
 Prometheus Operator provided by Red Hat The Prometheus Operator for Kubernetes provides easy monitoring definitions for	 Strimzi Kafka provided by Red Hat Run an Apache Kafka cluster, including Kafka Connect, ZooKeeper and more.	

Operator Development

- Operator can be implemented in any language
- Frameworks:
 - Operator Framework (Golang, Helm, Ansible)
<https://github.com/operator-framework>
 - Kubebuilder (Golang)
<https://github.com/kubernetes-sigs/kubebuilder>
 - Metacontroller (Language agnostic)
<https://metacontroller.app/>
 - jvm-operators (Java, Groovy, Kotlin, ...)
<https://github.com/jvm-operators>



Elastic
Scale

How to automatically
react to dynamically
changing resource
requirements.

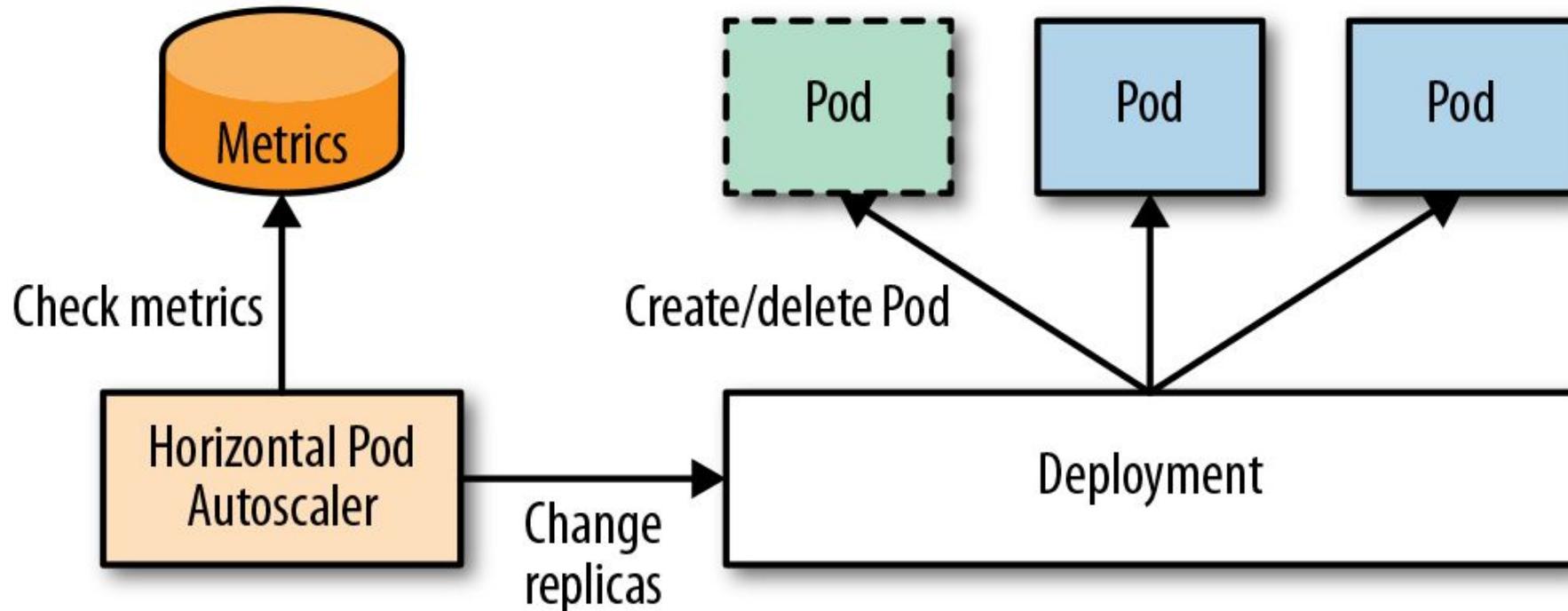
Scaling

- **Horizontal:** Changing replicas of a Pod
- **Vertical:** Changing resource constraints of containers in a single Pod
- **Cluster:** Adding new nodes to a cluster
- **Manual:** Changing scale parameters manually,. imperatively or declaratively
- **Automatic:** Change scaling parameters based on observed metrics

Scaling

- **Horizontal:** Changing replicas of a Pod
- **Vertical:** Changing resource constraints of containers in a single Pod
- **Cluster:** Adding new nodes to a cluster
- **Manual:** Changing scale parameters manually,. imperatively or declaratively
- **Automatic:** Change scaling parameters based on observed metrics

Horizontal Pod Autoscaler (HPA)



```
kubectl autoscale deployment random-generator --cpu-percent=50 --min=1 --max=5
```

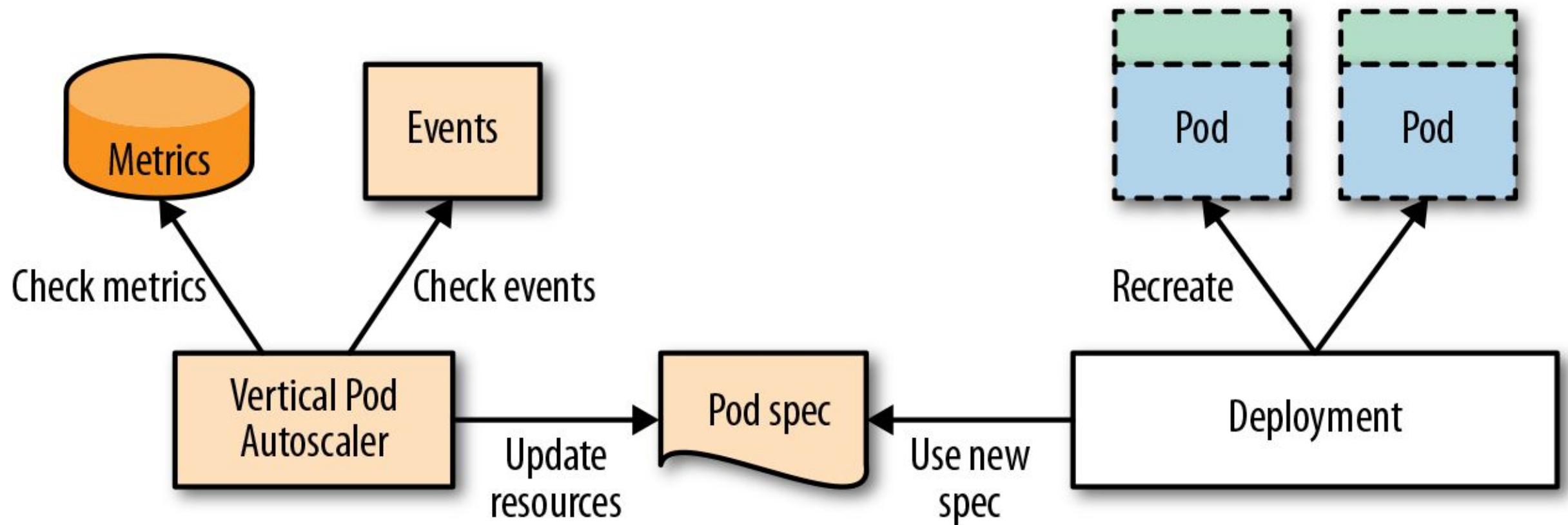
HorizontalPodAutoscaler

```
apiVersion: autoscaling/v2beta2
kind: HorizontalPodAutoscaler
metadata:
  name: random-generator
spec:
  minReplicas: 1
  maxReplicas: 5
  scaleTargetRef:
    apiVersion: extensions/v1beta1
    kind: Deployment
    name: random-generator
  metrics:
  - resource:
      name: cpu
      target:
        averageUtilization: 50
        type: Utilization
      type: Resource
```

HPA: Metrics & Challenges

- Metrics
 - **Standard Metrics** - CPU & Memory Pod data obtained from Kubernetes metrics server
 - **Custom Metrics** - Metrics delivered via an aggregated API server at the custom.metrics.k8s.io API path
 - **External Metrics** - Metrics obtained from outside the cluster
- Challenges
 - **Metric Selection** - Correlation between metric value and replica counts
 - **Preventing Thrashing** - Windowing to avoid scaling on temporary spikes
 - **Delayed Reaction** - Delay between cause and scaling reaction

Vertical Pod Autoscaler (VPA)



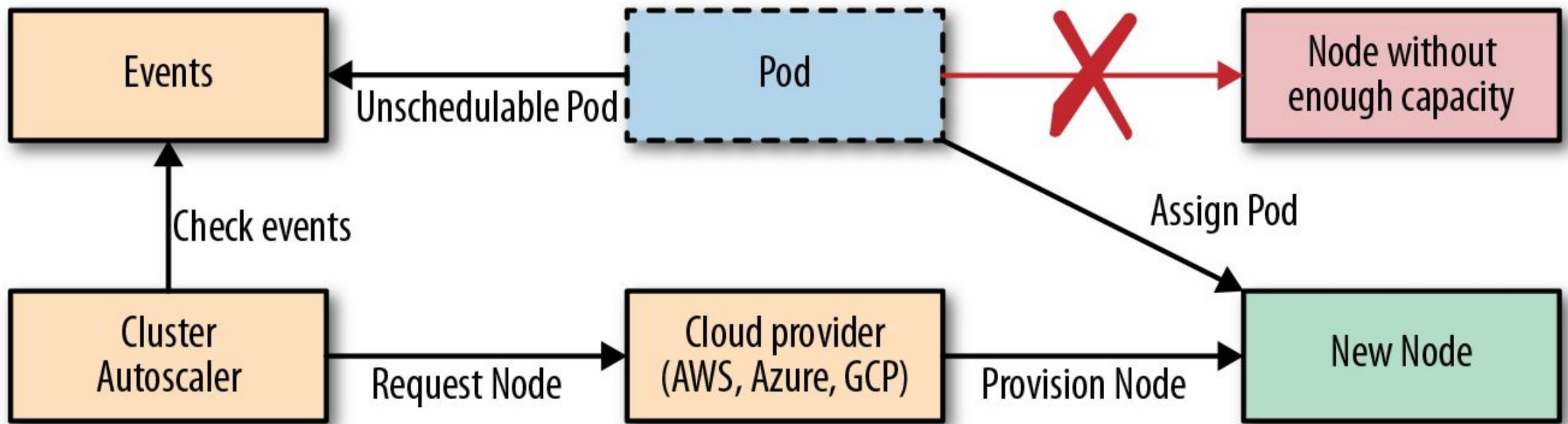
VerticalPodAutoscaler

```
apiVersion: poc.autoscaling.k8s.io/v1alpha1
kind: VerticalPodAutoscaler
metadata:
  name: random-generator-vpa
spec:
  selector:
    matchLabels:
      app: random-generator
  updatePolicy:
    updateMode: "Off"
```

VPA: Update Mode

- updateMode : Off
 - Recommendations are stored in the status : section of the VPA resource
 - No changes to the selected resources are performed
- updateMode : Initial
 - Recommendations are applied during creation of a Pod
 - Influences scheduling decision
- updateMode : Auto
 - Automatically restarts Pods with updated resources based on recommendation

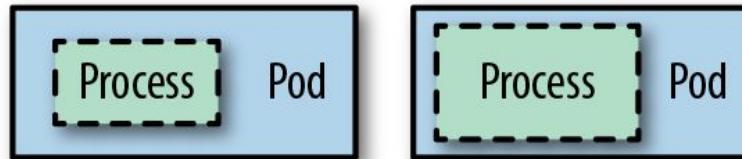
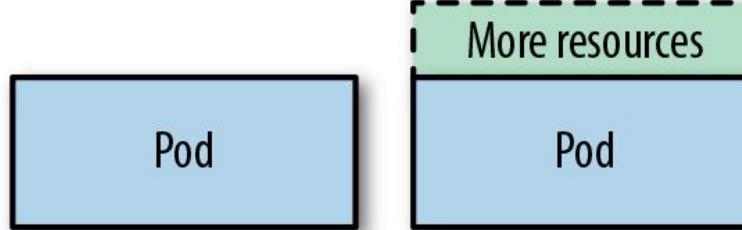
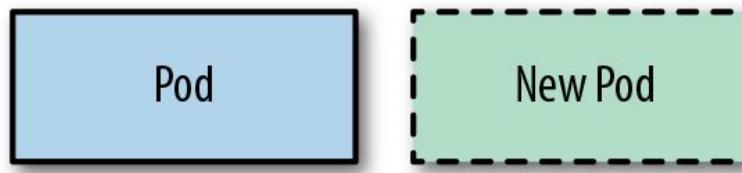
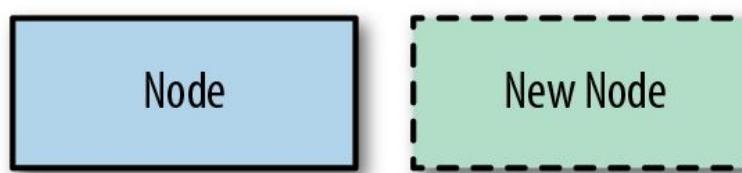
Cluster Autoscaler



Cluster Autoscaler

- Scale-Up
 - Adding a new node if a Pod is marked as *unschedulable* because of scarce resources.
 - Cluster API: Kubernetes API for dynamically managing node groups.
- Scale-Down
 - ... more than half of a nodes capacity is unused
 - ... all movable Pods can be placed on other nodes
 - ... no other reasons to prevent node deletion
 - ... no Pods that can not be moved

Elastic Scale

Technique	Action	Scaling Example
Application Tuning	Tune process (threads, heap, etc.)	
Vertical Pod Autoscaler	Increase/reduce container resources	
Horizontal Pod Autoscaler	Add/remove Pods	
Cluster Autoscaler	Add/remove Nodes	

Thank you



<https://k8spatterns.io>



@ro14nd



@bibryam



@k8spatterns

Picture Credits

<https://www.pexels.com/photo/brown-and-black-pattern-2158386/>

<https://pixabay.com/photos/ship-helm-sunset-cutter-coast-guard-759954/>

<https://unsplash.com/photos/yo01Z-9HQAw>

<https://www.pexels.com/photo/turned-on-light-crane-1117211/>

<https://www.pexels.com/photo/shallow-focus-photography-of-black-ship-1095814/> <https://unsplash.com/photos/UHDx3BHIFvY>

<https://pixabay.com/photos/containers-storage-rusted-rusty-1209079/>

<https://pixabay.com/photos/motocross-sidecar-race-motorsport-1045661/>

<https://www.pexels.com/photo/golden-gate-bridge-san-francisco-california-1141> https://unsplash.com/photos/ymf4_9Y9S_A

<https://www.pexels.com/photo/reflection-playstation-pad-gaming-18174/>

https://unsplash.com/photos/UJP_QpCKj7M

<https://www.pexels.com/photo/grayscale-photo-of-person-holding-chess-piece-1498958/>

<https://pixabay.com/photos/cans-manufacturing-business-2888650/>

<https://unsplash.com/photos/IRoX0shwjUQ>

<https://www.freeimages.com/photo/poppy-in-wheat-1344010>

<https://pixabay.com/photos/telescope-field-glass-spyglass-1966366/>

<https://unsplash.com/photos/UHDx3BHIFvY>

<https://pixabay.com/photos/bake-advent-christmas-cookie-1786926/>

<https://pixabay.com/photos/balloons-colors-party-celebration-1869790/>