



# ***Disaster Recovery Options running Apache Kafka On Kubernetes***

**Jennifer Snipes**

Staff Customer Success Technical Architect

**Rema Subramanian**

Customer Success Technical Architect

# Contents



**1. Resilient Kafka Architectures**

**2. Kafka & Kubernetes**

**3. Getting Started**

**4. Kubernetes Operator**

**5. Stretch Cluster on Kubernetes**

**6. Putting it to the Test**

**7. Wrapping it Up**

**8. Demo / Q & A**



# *Resilient Kafka Architectures*

# Know your RTO & RPO



I NOTICE IN YOUR  
DISASTER PREVENTION PLAN  
THE RECOVERY OBJECTIVE  
TIME IS 'WHENEVER'.



## RPO

RPO (Recovery Point Objective) is about **how much data you can afford to lose before it impacts business operations**. For example, for a banking system, 1 hour of data loss can be catastrophic as they operate live transactions.

## RTO

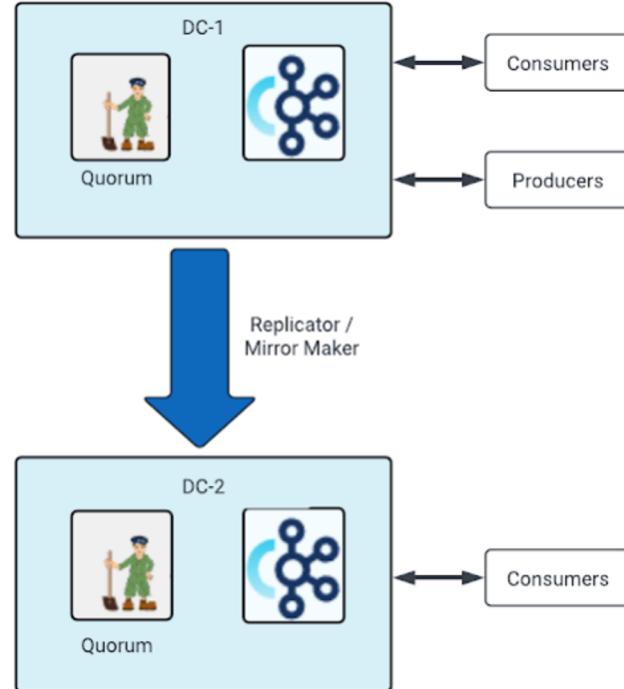
RTO (Recovery Time Objective) is **the timeframe within which an application and systems must be restored** after an outage.

# Resilient Kafka Architectures



## Active DC-1 / Passive DC-2

- Two independent clusters in different Data Centers / Regions
- Producers only in one Data Center
- Consumers in both Data Centers
- Multi-cloud / Multi-region
- One way Replication
- Asynchronous Replication
- RPO >0, RTO >0

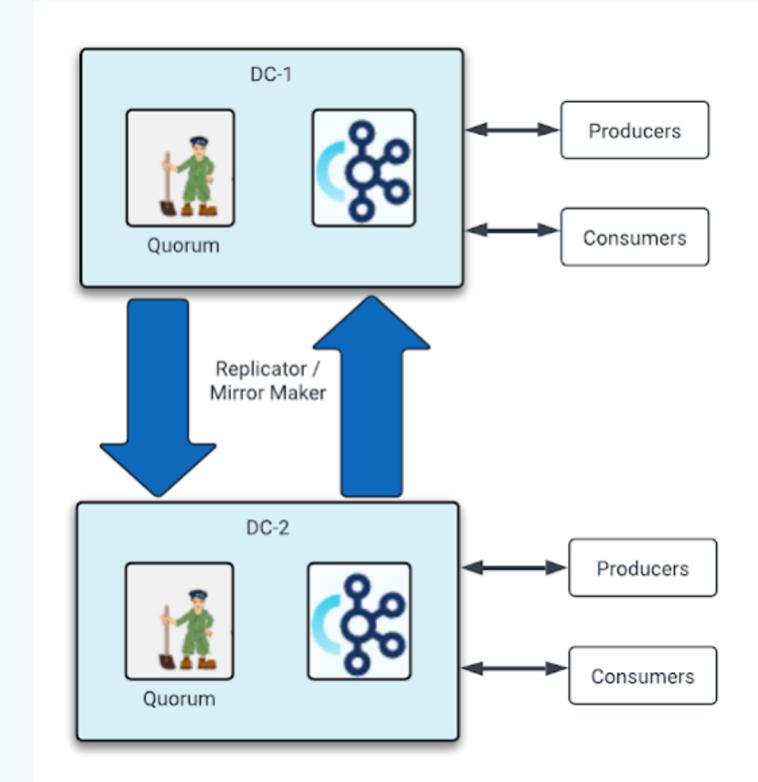


# Resilient Kafka Architectures



## Active DC-1 / Active DC-2

- Two independent clusters in different Data Centers / Regions
- Producers in both Data Centers
- Consumers in both Data Centers
- Multi-region / Multi-cloud
- Bi-Directional Replication with Provenance Headers
- Asynchronous Replication
- RPO >0, RTO >0

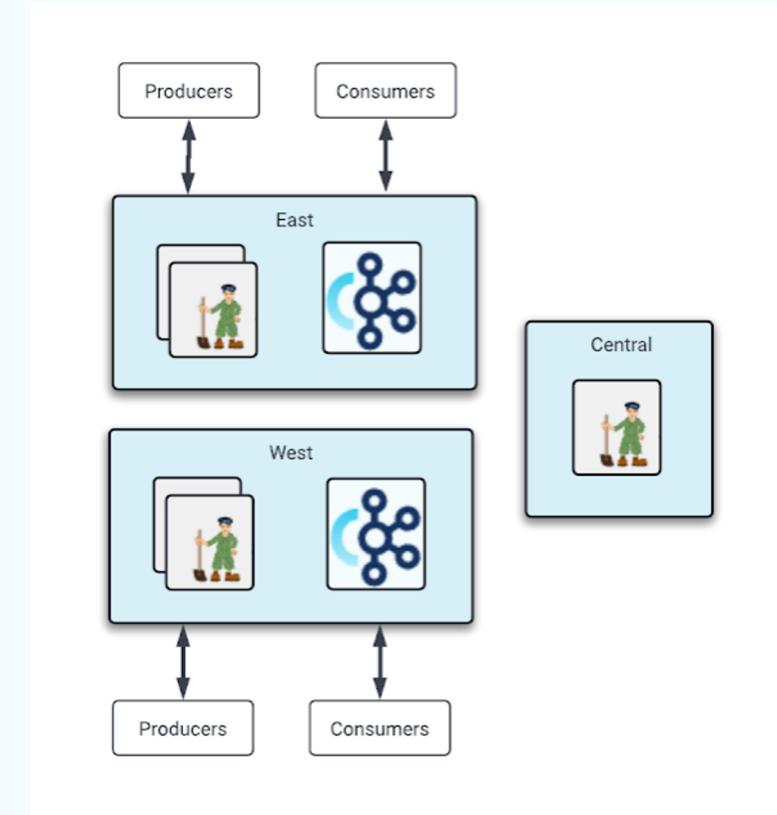


# Resilient Kafka Architectures



## Stretch Cluster

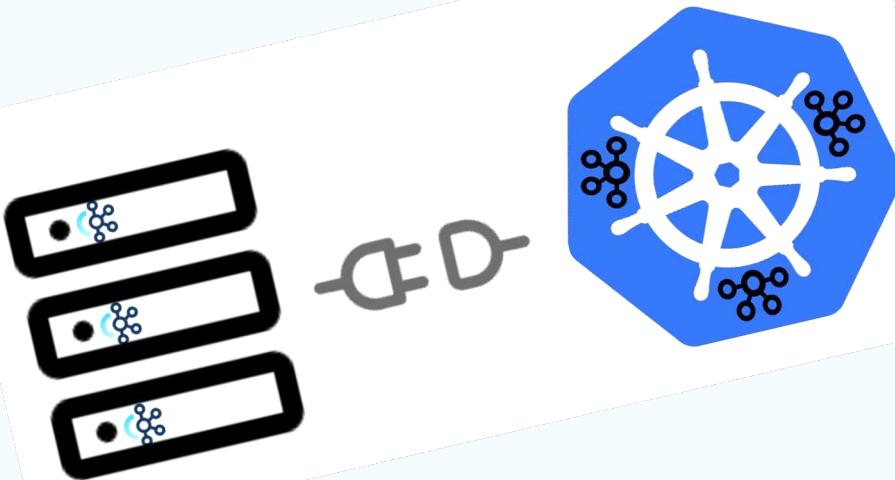
- Single Cluster stretched across different Data Centers
- Producers write transparently across Data Centers
- Consumers in all Data Centers
- RPO = 0, RTO near 0
- Synchronous Replication native to cluster
- Asynchronous Replication with Observers\*
- Replica & Observer placement defines Active-Active vs Active-Passive\*
- Auto Observer Promotion\*
- Multi-region\*





# *Kafka & Kubernetes*

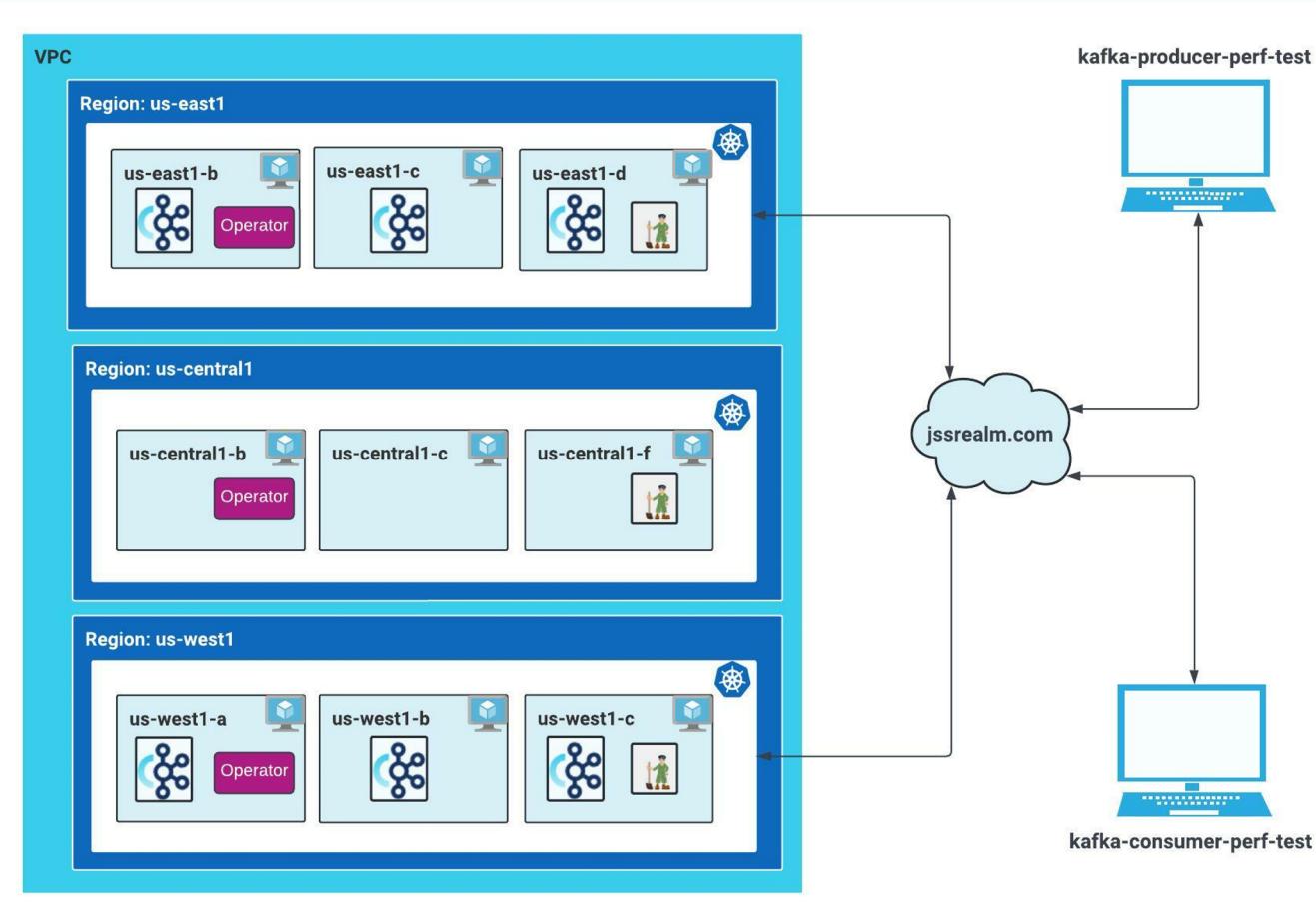
# Kafka Deployment Arenas



## Traditional vs K8s

- Broker
  - Hostnames /IPs
  - Placement across DCs
  - Communication across DCs
  - Failure
- Broker and ZK co-location
- Multi-tenancy

# Stretch Cluster on Kubernetes



# *Built-in Disaster Recovery on Kubernetes*





# *Getting Started*

# *Building the K8s Cluster*



1

## GCP VPC Native Cluster

- Alias IP address range for nodes, pods and services
- Requires non-overlapping CIDR ranges



2

## GKE Cluster

- Separately Managed Node Pool
- Node machine type
- Configurable number of nodes distributed across AZs
- Distinct namespace per cluster



3

## Networking

- VPC Native cluster installs routing
- Firewall rules
  - allow tcp between k8s clusters
  - allow access to 2181, 2888, 3888, 9092, 7778, 3000

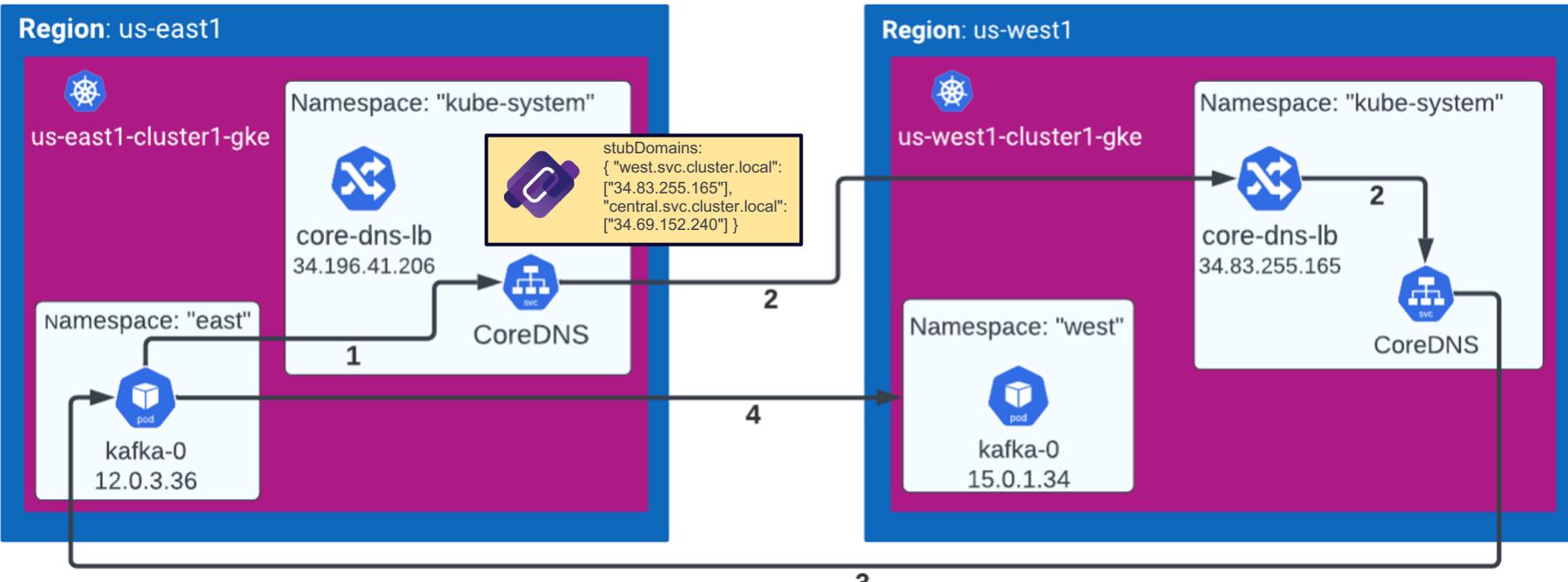


4

## StorageClass

- provisioner: kubernetes.io/gce-pd
- allowVolumeExpansion : true
- type: pd-ssd
- fstype: ext4
- reclaimPolicy: Retain
- volumeBindingMode: WaitForFirstConsumer

# Networking between Kubernetes Clusters





# *Kubernetes Operator*

# Operator



1

## CRDs

- Define various application components
- Medium to tie to kafka server.properties



2

## Controller

- CRD's behavior
- Reconciliation loop



3

## Services

- Headless Service
- Expose individual pods as external services
- Bootstrap LB service to get metadata



4

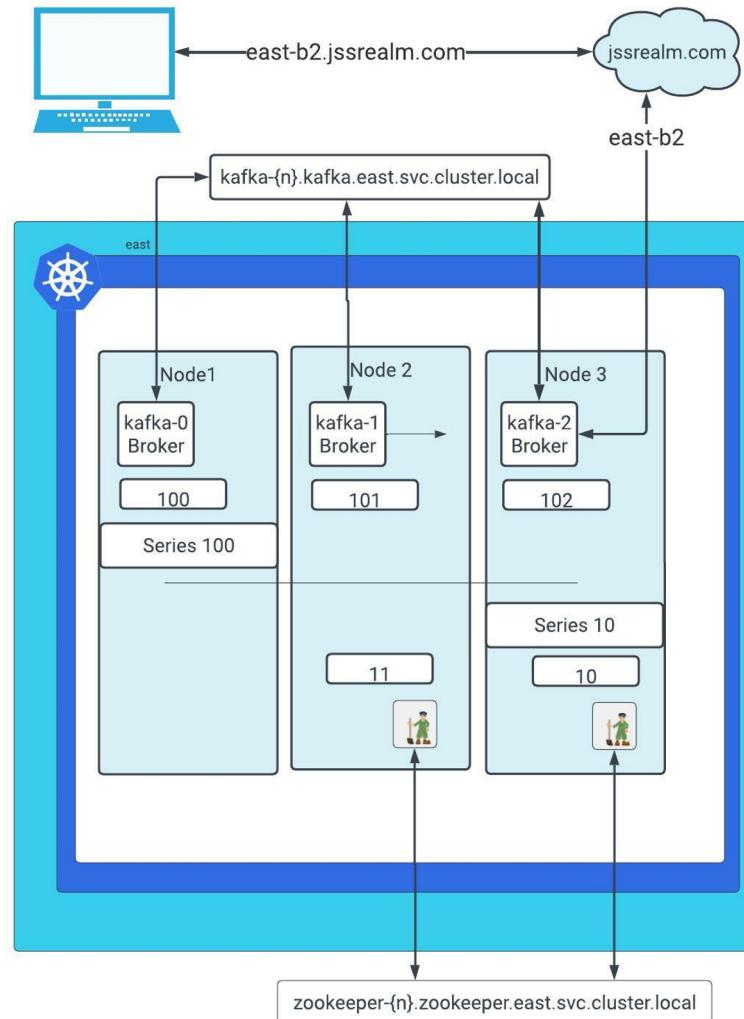
## StatefulSets

- PVC claim

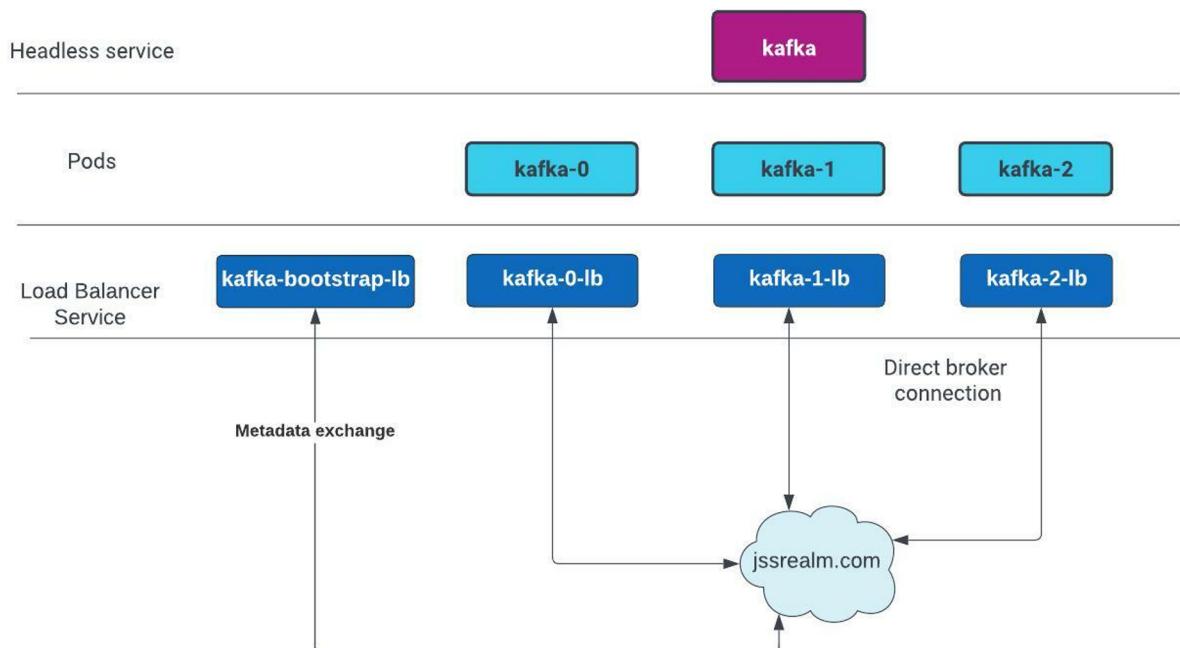
# Pod Accessibility

## Identifying the kafka pods

- Unique broker IDs
- Internally, each pod resolves `kafka-{n}.kafka.east.svc.cluster.local`
- Externally, broker prefixes to map to pod ordinals  
`{region}-b{n}.{domain}`



# Kafka CRD Services Galore





# *Stretch Cluster on Kubernetes*

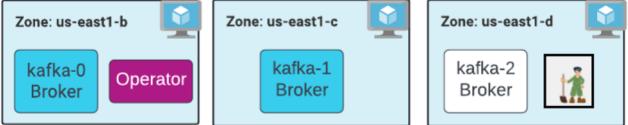
# Multi Region (Stretch) Cluster



## VPC-native cluster: us-cluster-vpc

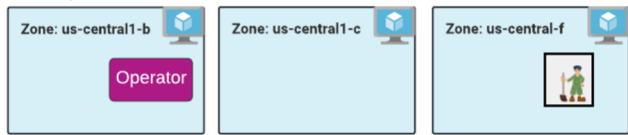
### Region: us-east1

Subnet: us-east1-cluster-subnet    Cluster: us-east1-cluster1-gke    Namespace: east  
CIDR ranges : primary(node) - 11.0.0.0/16; secondary (pod) - 12.0.0.0/16; secondary (service) - 13.0.0.0/20;



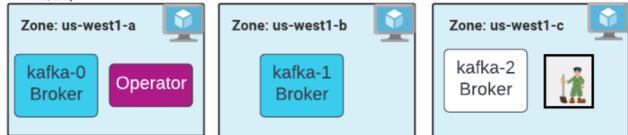
### Region: us-central1

Subnet: us-central1-cluster-subnet    Cluster: us-central1-cluster1-gke    Namespace: central  
CIDR ranges : primary(node) - 17.0.0.0/16; secondary (pod) - 18.0.0.0/16; secondary (service) - 19.0.0.0/20;



### Region: us-west1

Subnet: us-west1-cluster-subnet    Cluster: us-west1-cluster1-gke    Namespace: west  
CIDR ranges : primary(node) - 14.0.0.0/16; secondary (pod) - 15.0.0.0/16; secondary (service) - 16.0.0.0/20;



## East Cluster

- Single Kubernetes Cluster (us-east1-cluster1-gke)
- 3+ Brokers
- Single Zookeeper

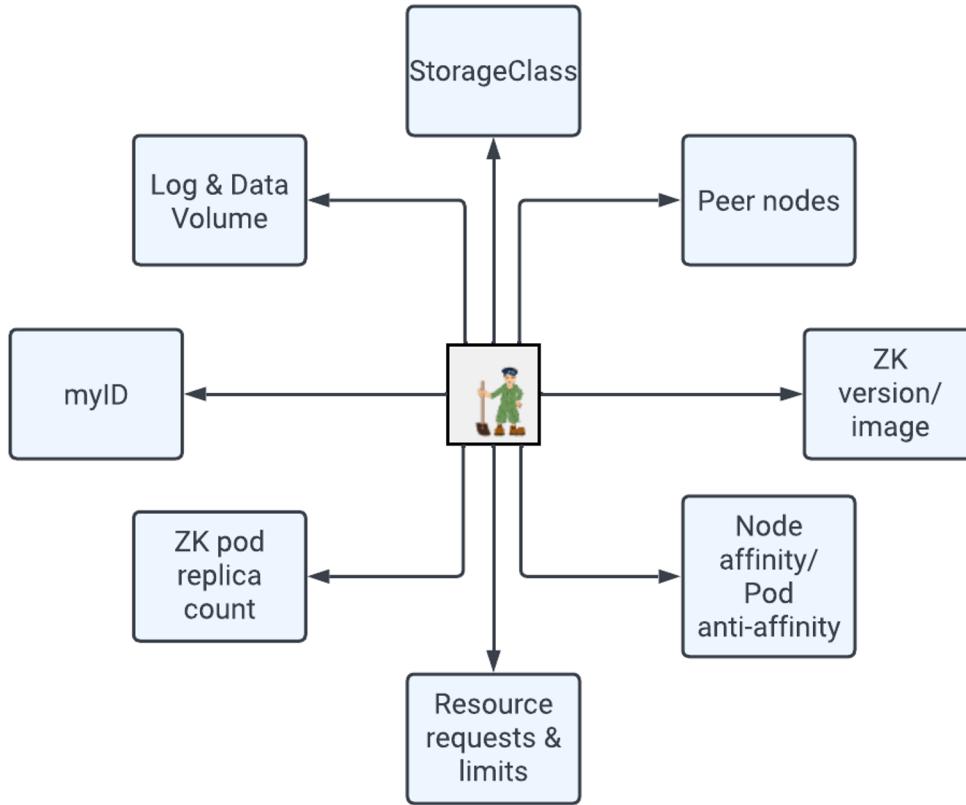
## Central Cluster

- Single Kubernetes Cluster (us-central1-cluster1-gke)
- Single Zookeeper

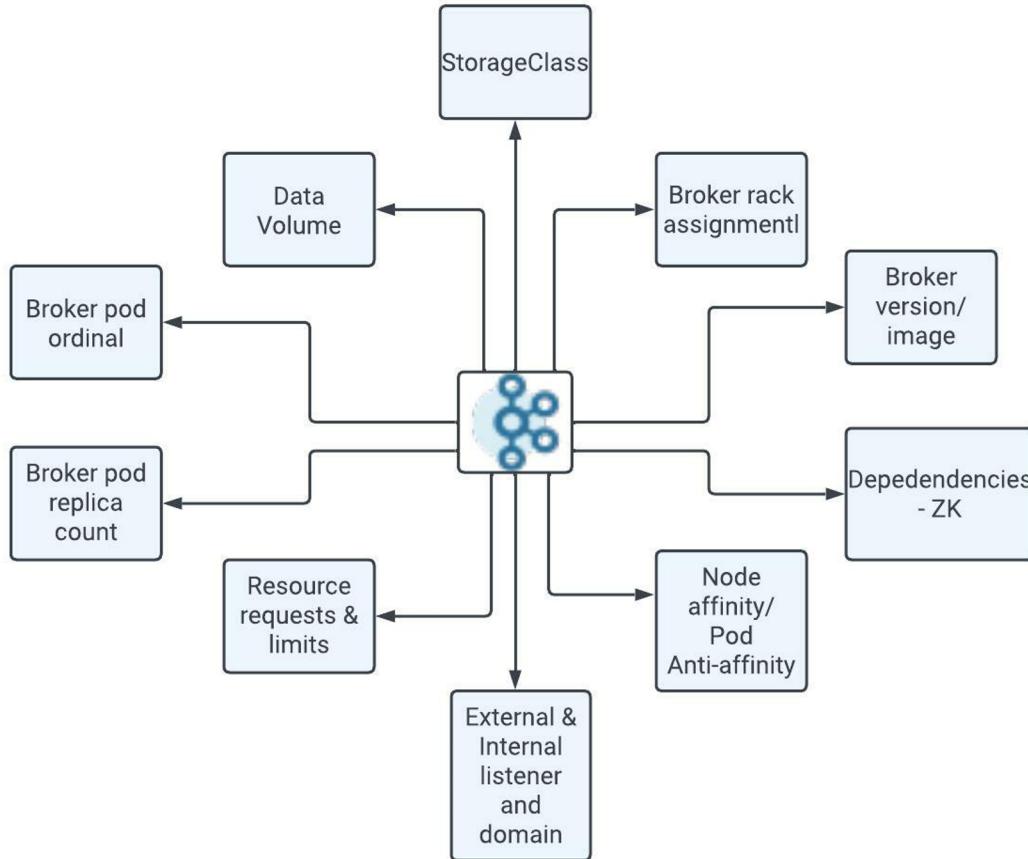
## West Cluster

- Single Kubernetes Cluster (us-west1-cluster1-gke)
- 3+ Brokers
- Single Zookeeper

# ZooKeeper Configuration



# Kafka Configuration



# Topic Replica Placement\*



```
1  [
2      "version": 2,
3      "replicas": [
4          {
5              "count": 1,
6              "constraints": {
7                  "rack": "us-east1-b"
8              },
9          },
10         {
11             "count": 1,
12             "constraints": {
13                 "rack": "us-east1-c"
14             },
15         },
16         {
17             "count": 1,
18             "constraints": {
19                 "rack": "us-west1-a"
20             },
21         },
22         {
23             "count": 1,
24             "constraints": {
25                 "rack": "us-west1-b"
26             },
27         },
28     ],
29     "observers": [
30         {
31             "count": 1,
32             "constraints": {
33                 "rack": "us-east1-d"
34             },
35         },
36         {
37             "count": 1,
38             "constraints": {
39                 "rack": "us-west1-c"
40             },
41         },
42     ],
43     "observerPromotionPolicy": "under-replicated"
44 ]
45 }
```

- **Broker Rack Awareness**
- **Synchronous Replicas**
- **Asynchronous Observers\***
- **Observer Promotion Policy\***



## *Putting it to the Test*

# Testing



## Pod Kill

- Deleted a pod

## Node Fault

- Node VM down
- Auto-scaler off & Node VM down

## Pod Failure

- Introduced pod failures with chaos-mesh

## ZK Quorum Failure

- 2 ZK nodes down - disrupt quorum

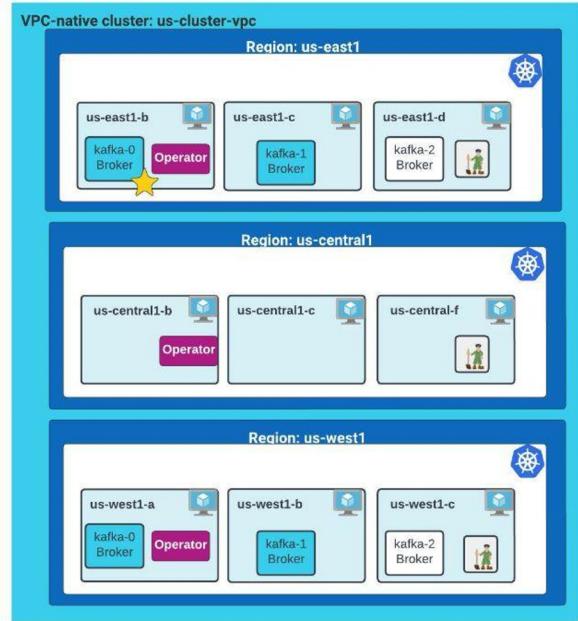
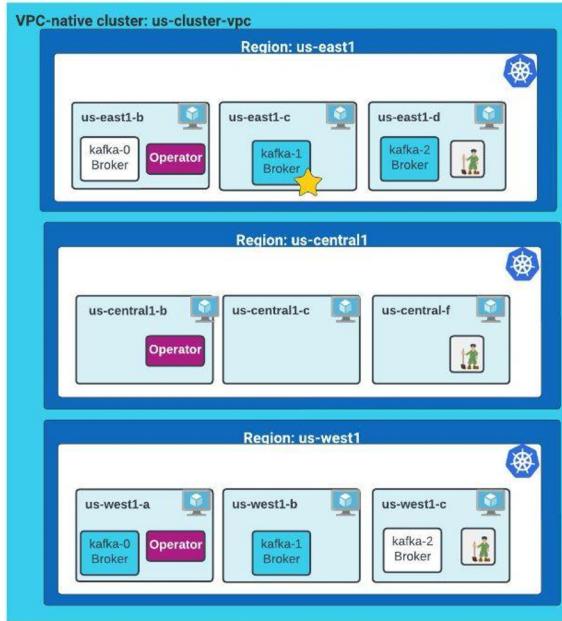
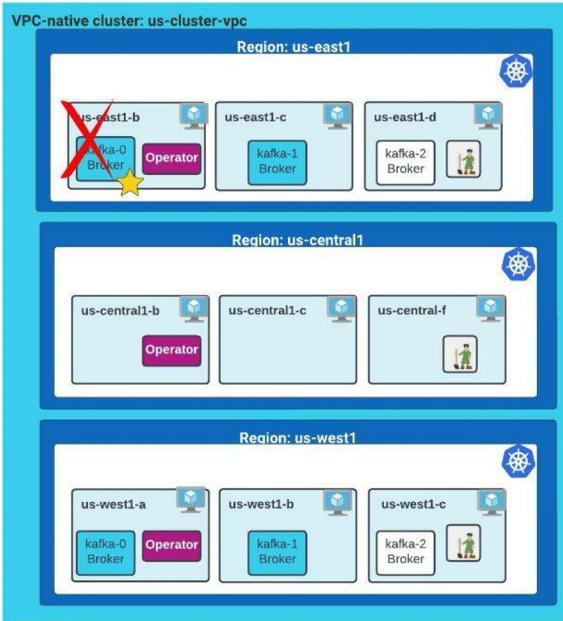
## Region Down

- Producers don't stop
- Controller broker and topic leader broker move to west
- ZK west is accessible for write

## Network Fault

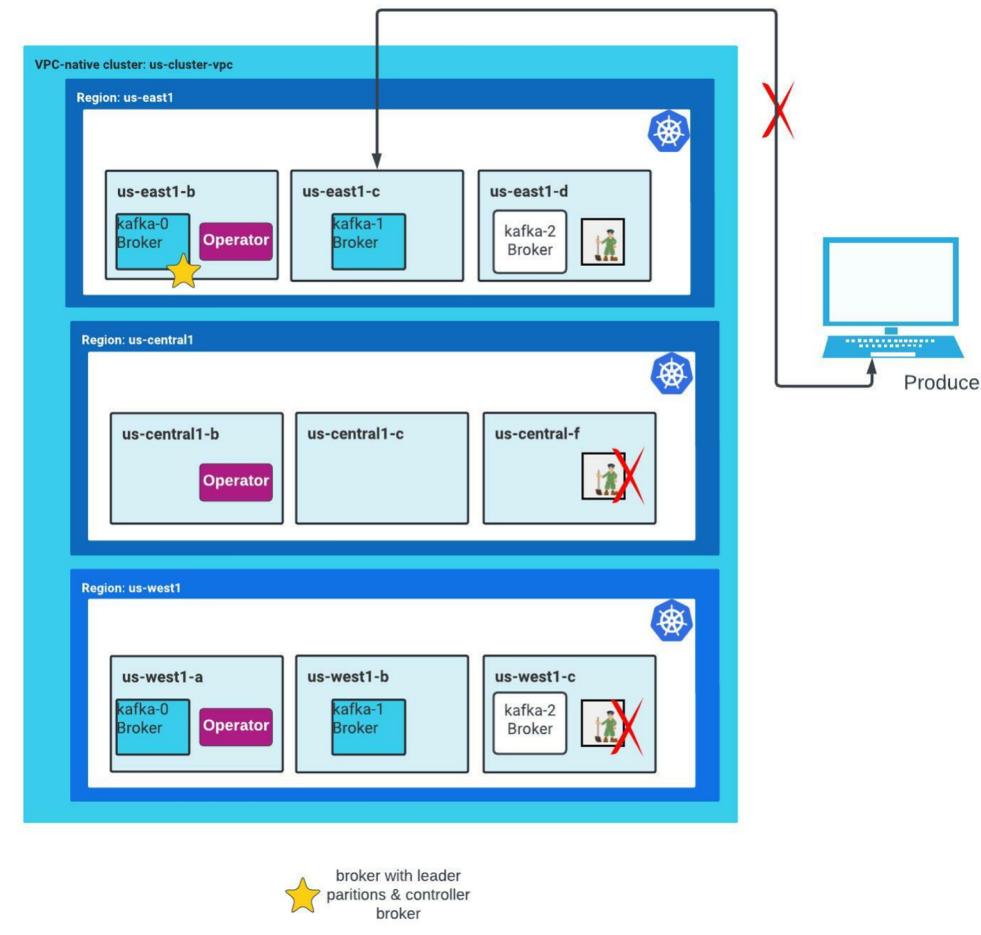
- Edited kube-dns stub-domain

# Pod Failure/Node Fault

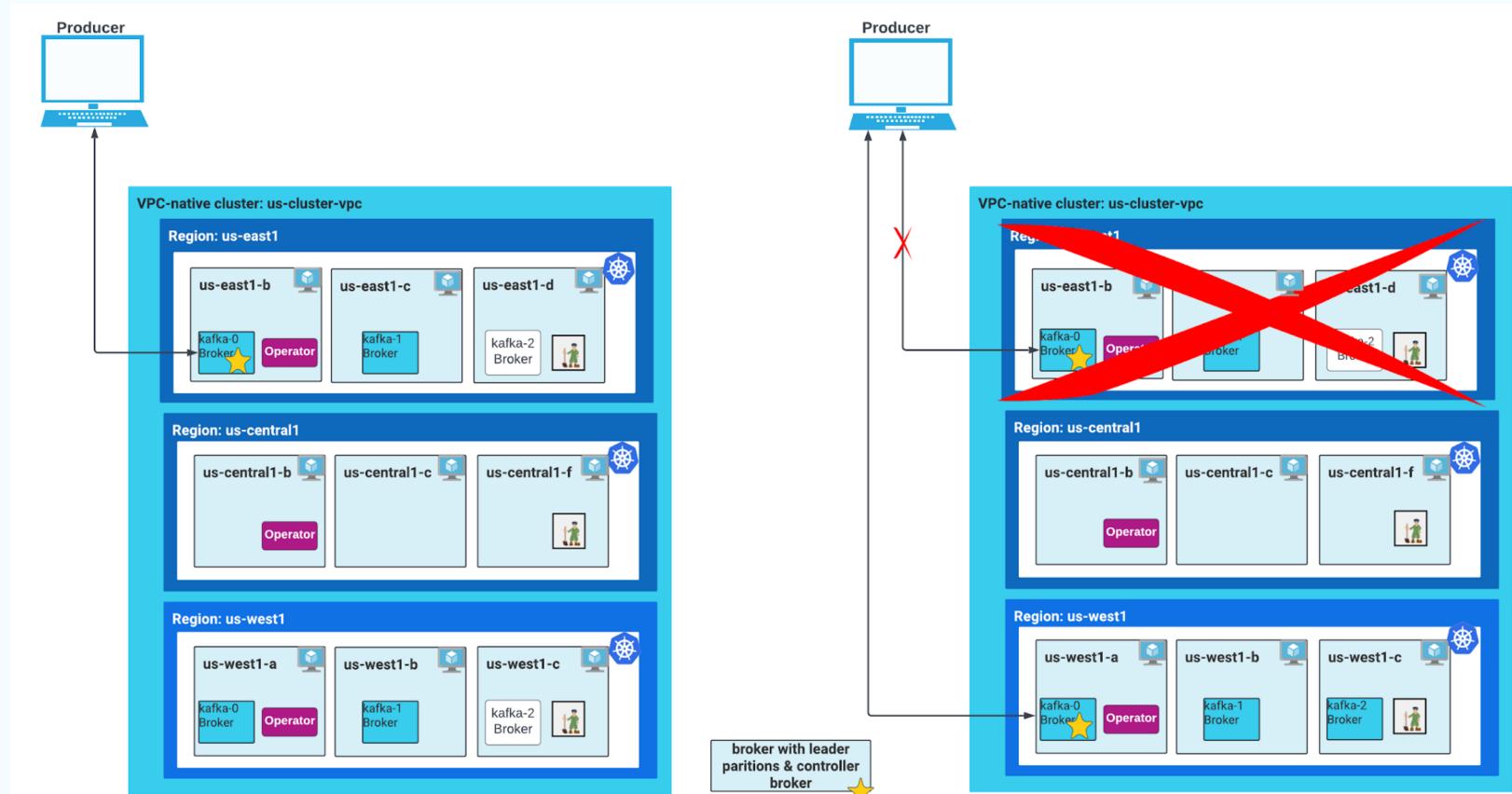


★ broker with leader partitions &  
controller broker

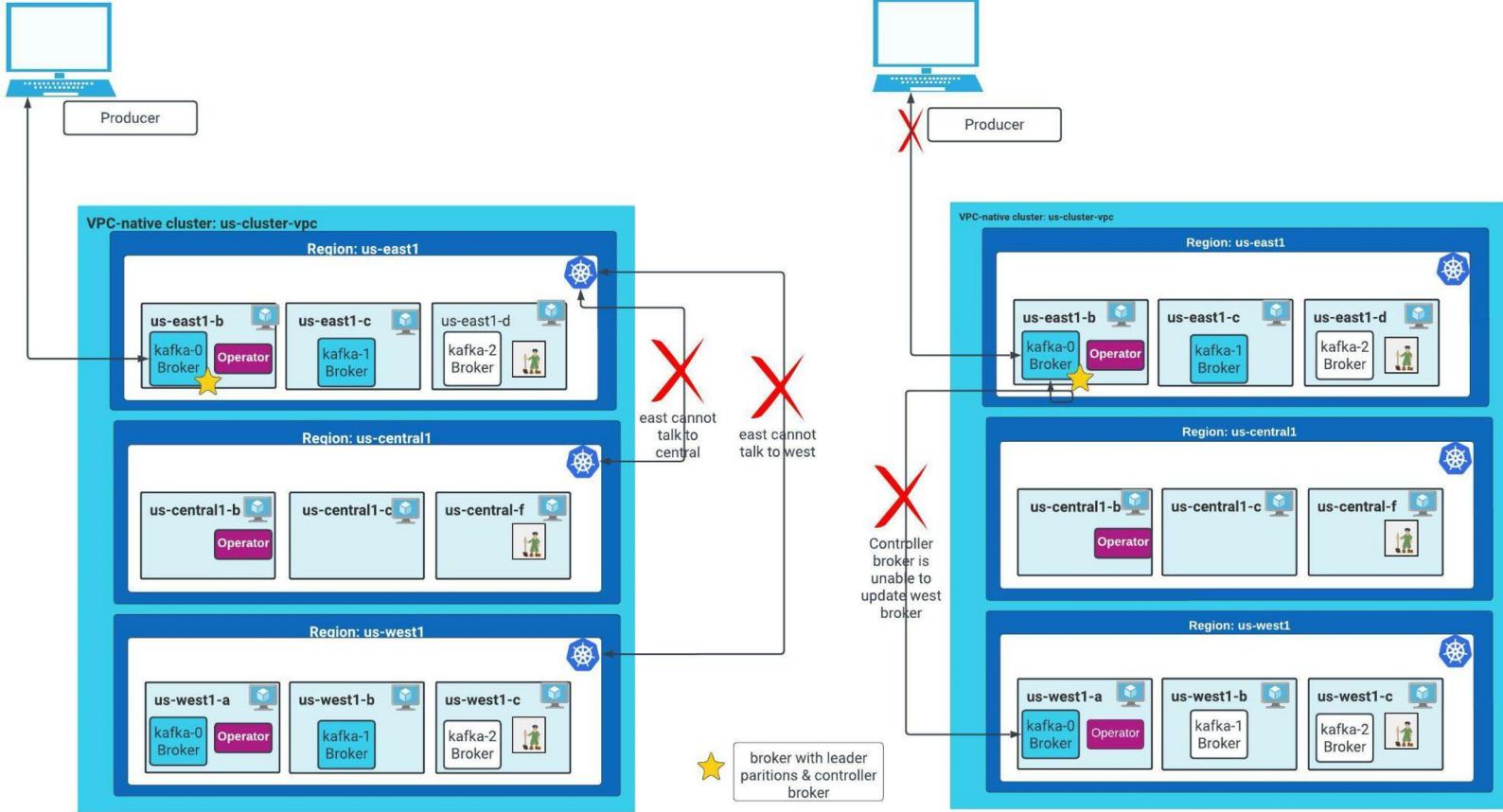
# ZK Quorum Failure



# Region Down



# Network fault



# Watching the Metrics





# *Wrapping it Up*

# Best Practices - Node to Pod Ratio



## Quantity

- Use eventsizer.io to derive number of broker pods
- Adjust count to balance across AZs
- Each AZ is a rack

## Capacity

- Eventsizer.io output will derive CPU and memory per pod
- Set CPU and memory limits and requests

## Nodes

- Memory optimized node type
- Evaluate capacity of node based on how average load
- Enable auto scaling with average and peak range



Pod

Pod

Nodes

# Best Practices



- Use Confluent for Kubernetes
- Choose the right storage that guarantees reliability, efficiency and speed e.g. SSD
  - Refer to 'Building the K8s Cluster' slide for other storage best practices
- Run health/liveness checks on:
  - Individual pod and bootstrap LBs
  - kube-system LBs
  - kube-dns service
- Use node affinity/pod anti-affinity to strategically place broker and ZK pods across AZs
- Use automation so infrastructure and CRDs can be deployed multi-region across all of your environments
- Follow best practices for tuning tcp socket buffers, replica fetcher, and clients for optimal performance with stretch clusters
- Minimum Durability Configuration
  - 2 Replicas and 2 Observers in each region
  - min.ISR=3
- Monitor everything!



# óotcha

- Requires separate IP CIDR ranges
- CoreDNS exposed externally
- Restricted to single K8s implementation
- CRDs may get stuck if the finalizer logic in Operator is not finishing
- Manually restart stateful set if pods are erroring - known issue
- GCP VPCs are global, subnets are regional

# *Achieving your Desired Resilience*



## Active-Passive

- RTO > 0, RPO > 0
- Replicas in one Region
- Observers in another Region
- Under-replicated AOP



## Active-Active

- RTO ~ 0, RPO = 0
- 2 Replicas in each region
- 2 Observers in each region
- Under-replicated AOP



# *Demo and Q&A*

# References



1. [https://assets.confluent.io/m/69c5ce7aff462f44/original/20180619-WP-Recommendations\\_for\\_Deploying\\_Apache\\_Kafka\\_on\\_Kubernetes.pdf](https://assets.confluent.io/m/69c5ce7aff462f44/original/20180619-WP-Recommendations_for_Deploying_Apache_Kafka_on_Kubernetes.pdf)
2. <https://cloud.google.com/kubernetes-engine/docs/how-to/alias-ips>
3. <https://learn.hashicorp.com/tutorials/terraform/gke>
4. <https://chaos-mesh.org/docs/>
5. <https://docs.confluent.io/operator/current/overview.html>
6. <https://docs.confluent.io/platform/current/multi-dc-deployments/multi-region.html>
7. <https://www.confluent.io/en-gb/events/kafka-summit-americas-2021/a-tale-of-2-n-data-centers-tuning-apache-kafka-clusters-to-combat-latency/>



*Your Apache Kafka®  
journey begins here*

[developer.confluent.io](https://developer.confluent.io)



CONFLUENT