



KubeCon



CloudNativeCon

North America 2023

Patterns of Multi-Cluster Kubernetes



Dan McKEAN

Sr Product Manager, MongoDB



George HANTZARAS

Engineering Director, MongoDB



Mircea COSBUC

Sr Engineer, MongoDB

Agenda

- Motivations for Multi-Cluster
- Architectures
- Network Centric Patterns
- Cluster Inventory
- Workload Distribution
- Networking
- Multi-Cluster Operators and Controllers
- Future of Multi-cluster

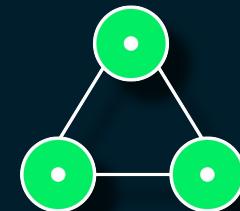
Motivations for Multi-Cluster



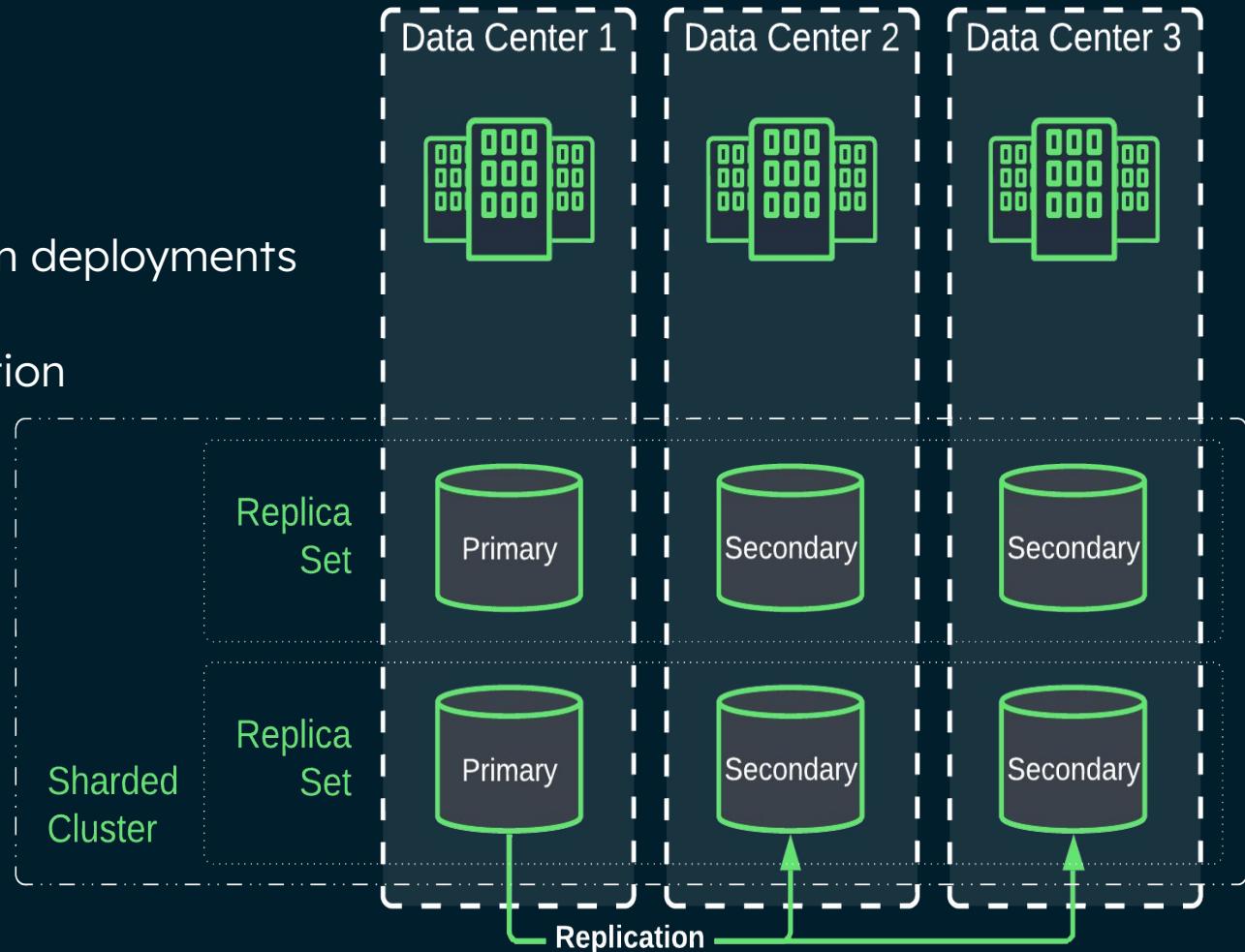


Motivations for Multi-Cluster Deployments

- Improved Performance & Reduced Latency
- Compliance with Regulations
- High Availability, Disaster Recovery, and Data Redundancy



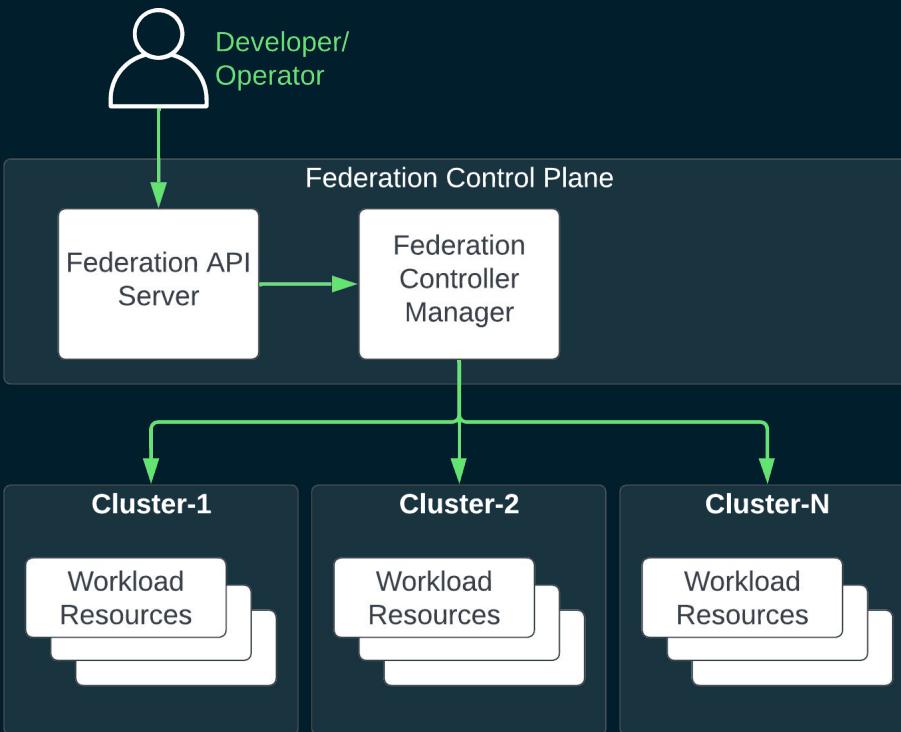
- Multi AZ, multi region deployments
- Cross region replication
- Segmenting data by location with multi-region shards
- Read-only



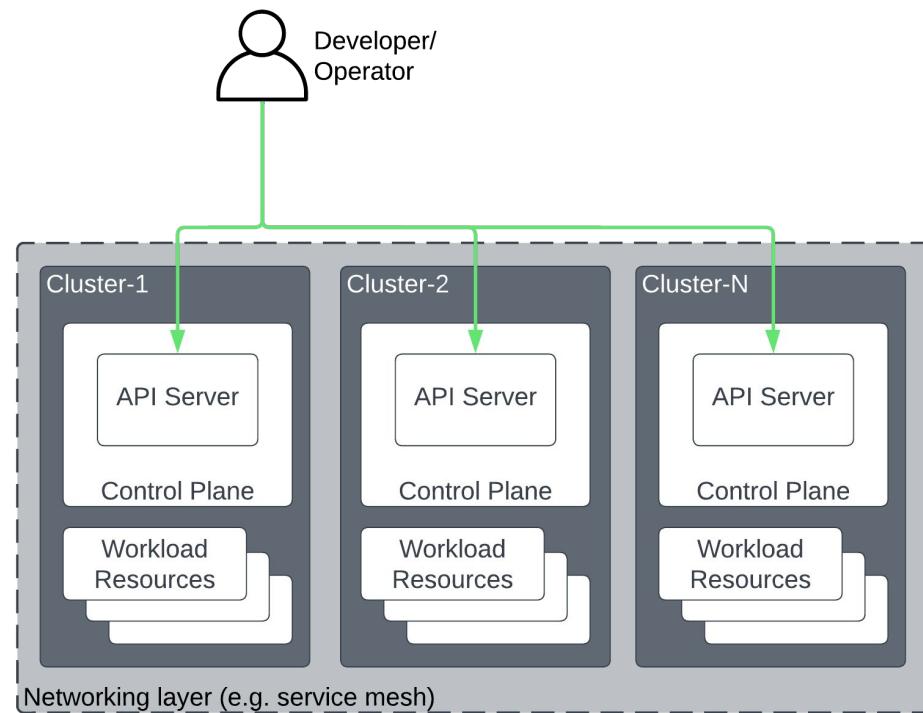
Architectures



Kubernetes-Centric (Cluster Federation)



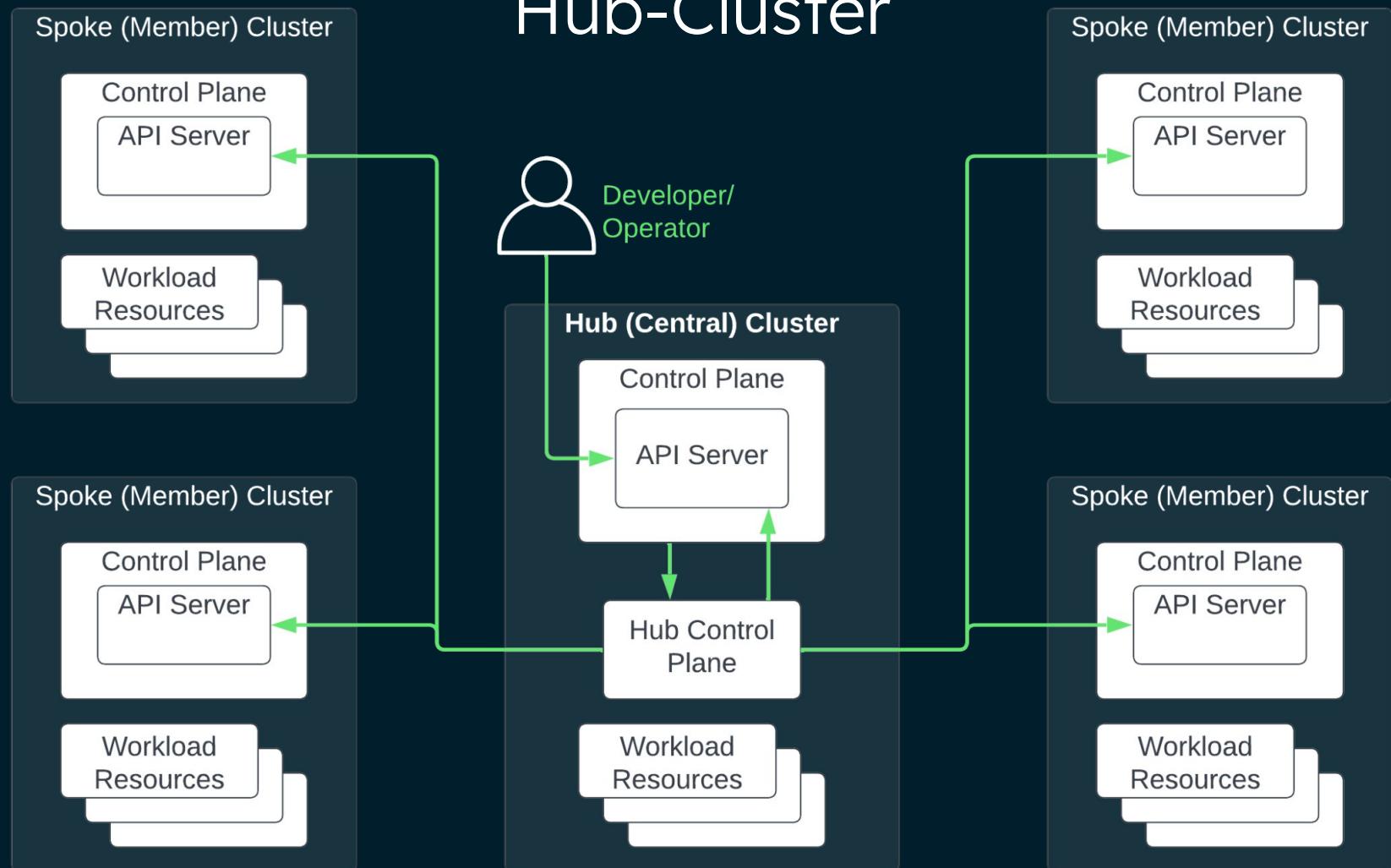
Network-Centric (Workload distribution)



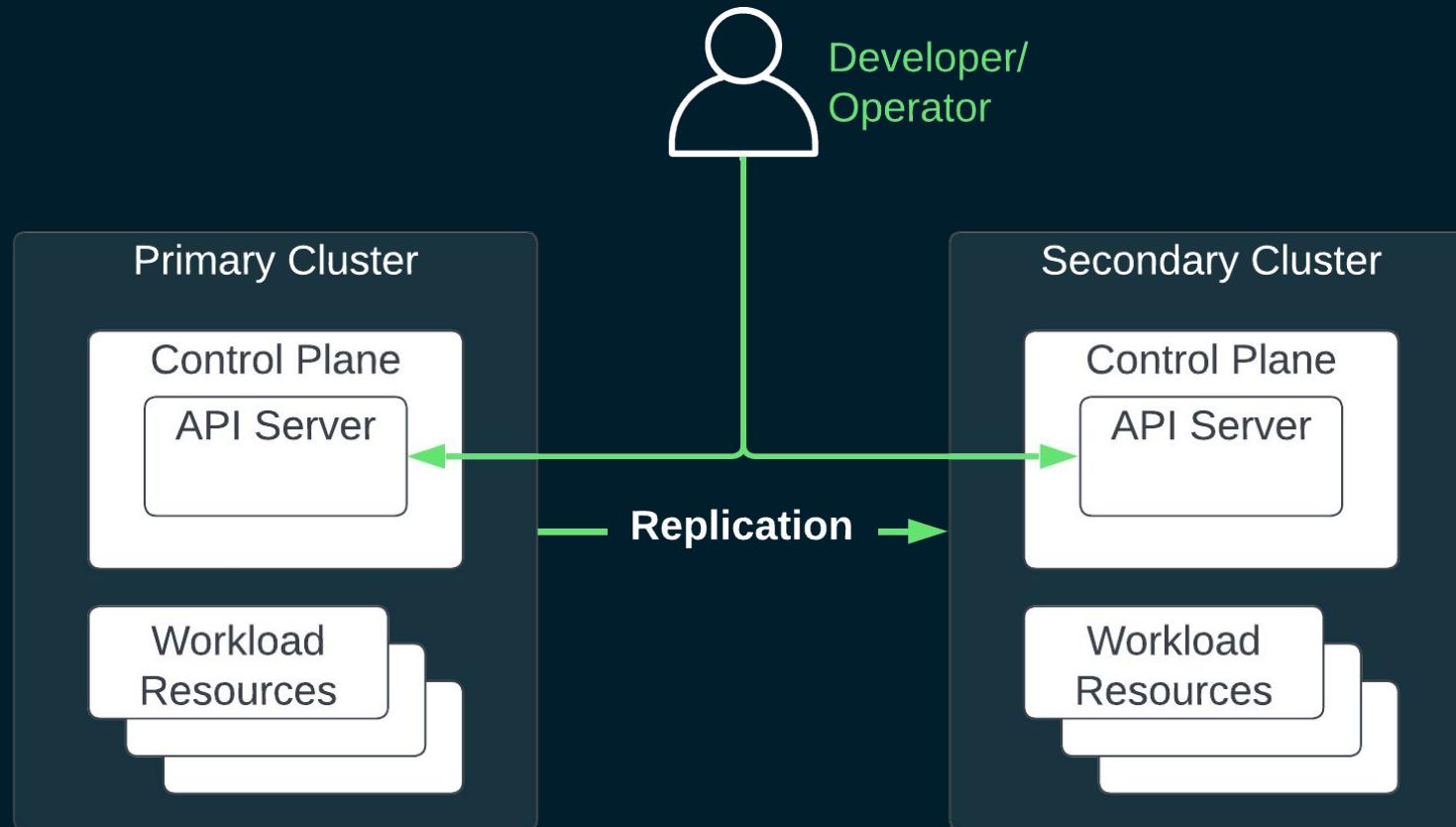
Network-Centric Patterns



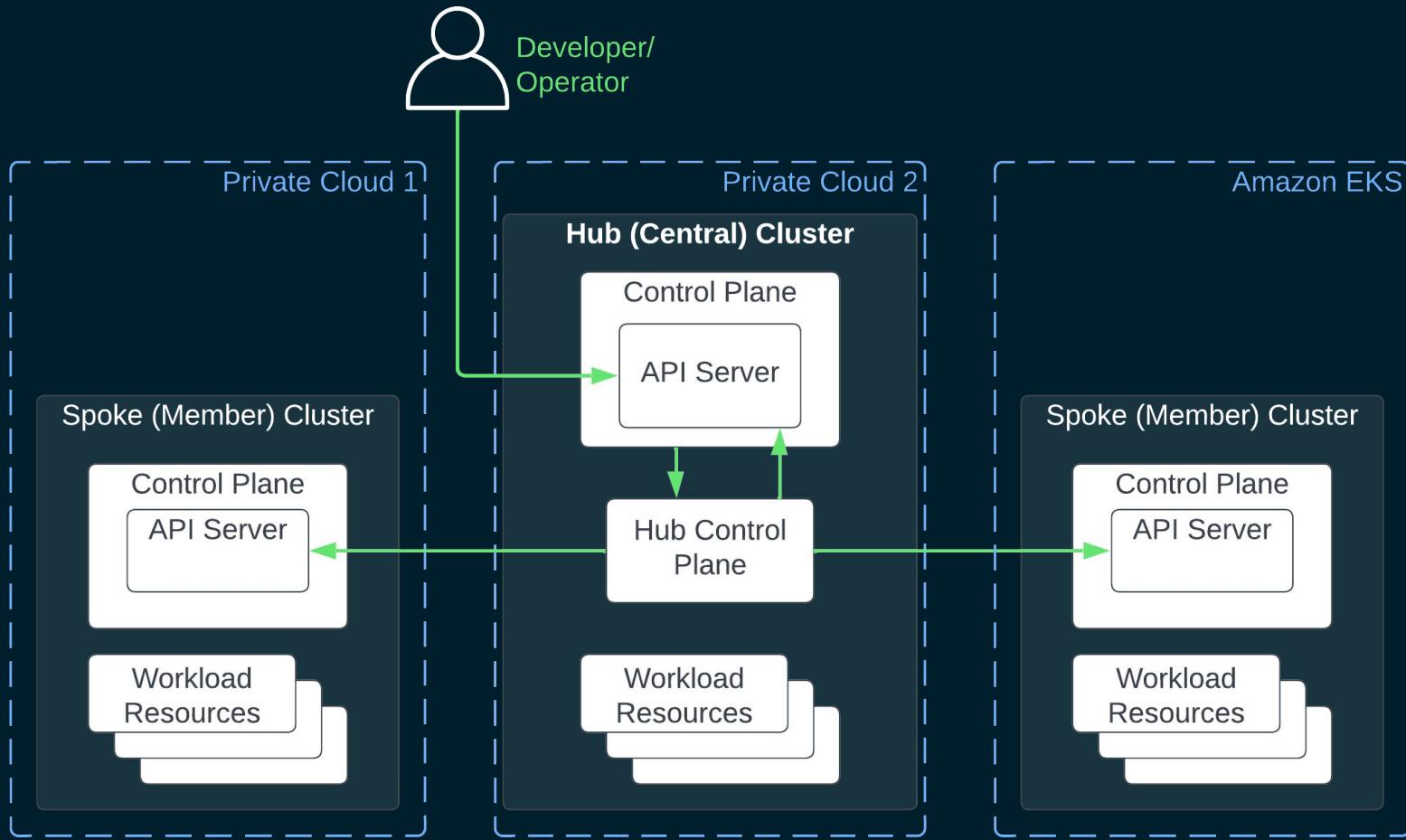
Hub-Cluster



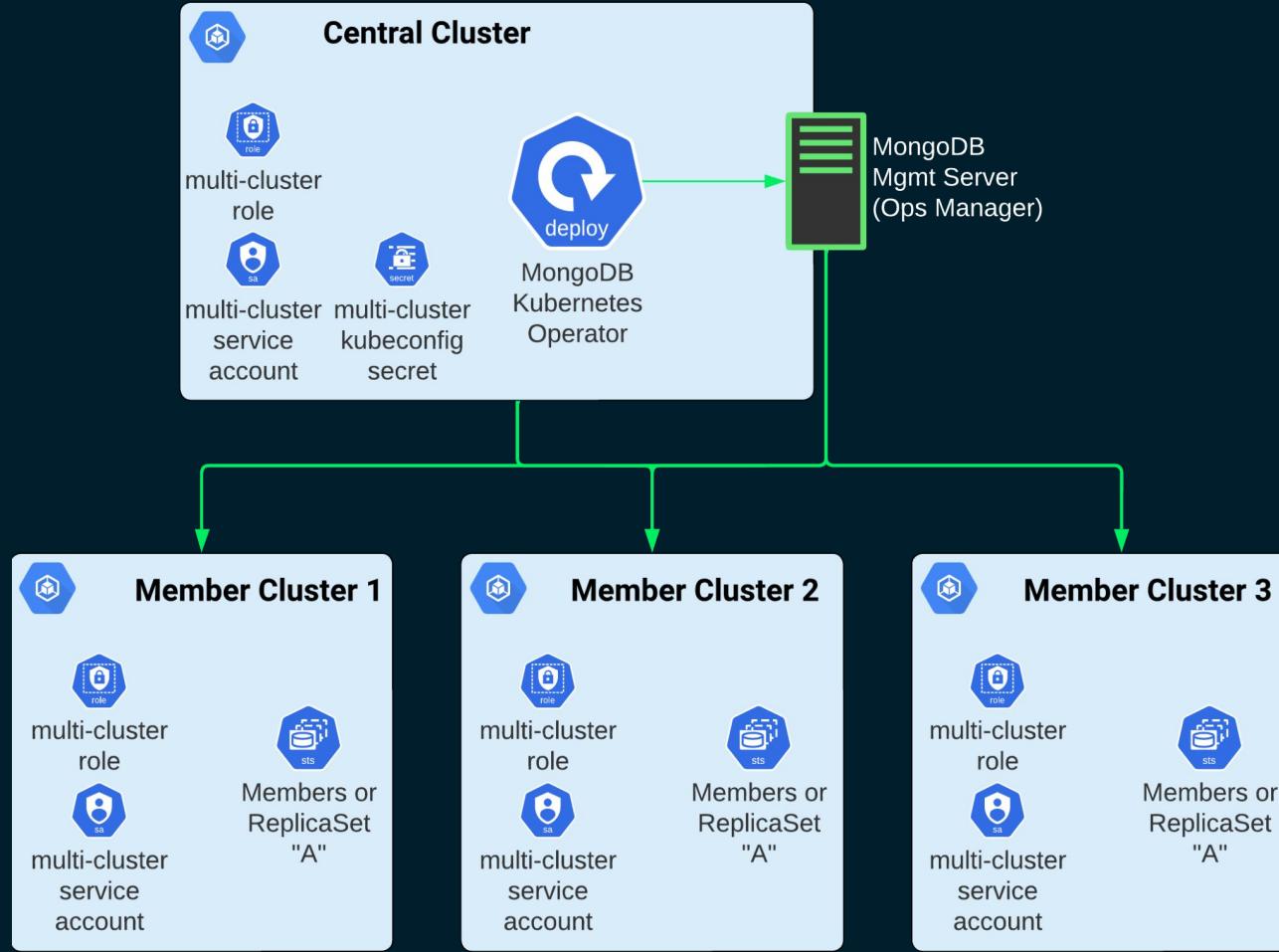
Replicated Application Architecture



Multi-cloud



MongoDB Architecture



Cluster Inventory



Cluster Inventory

Requirements



Cluster
credentials



Cluster health
checks



Cluster resource
tracking



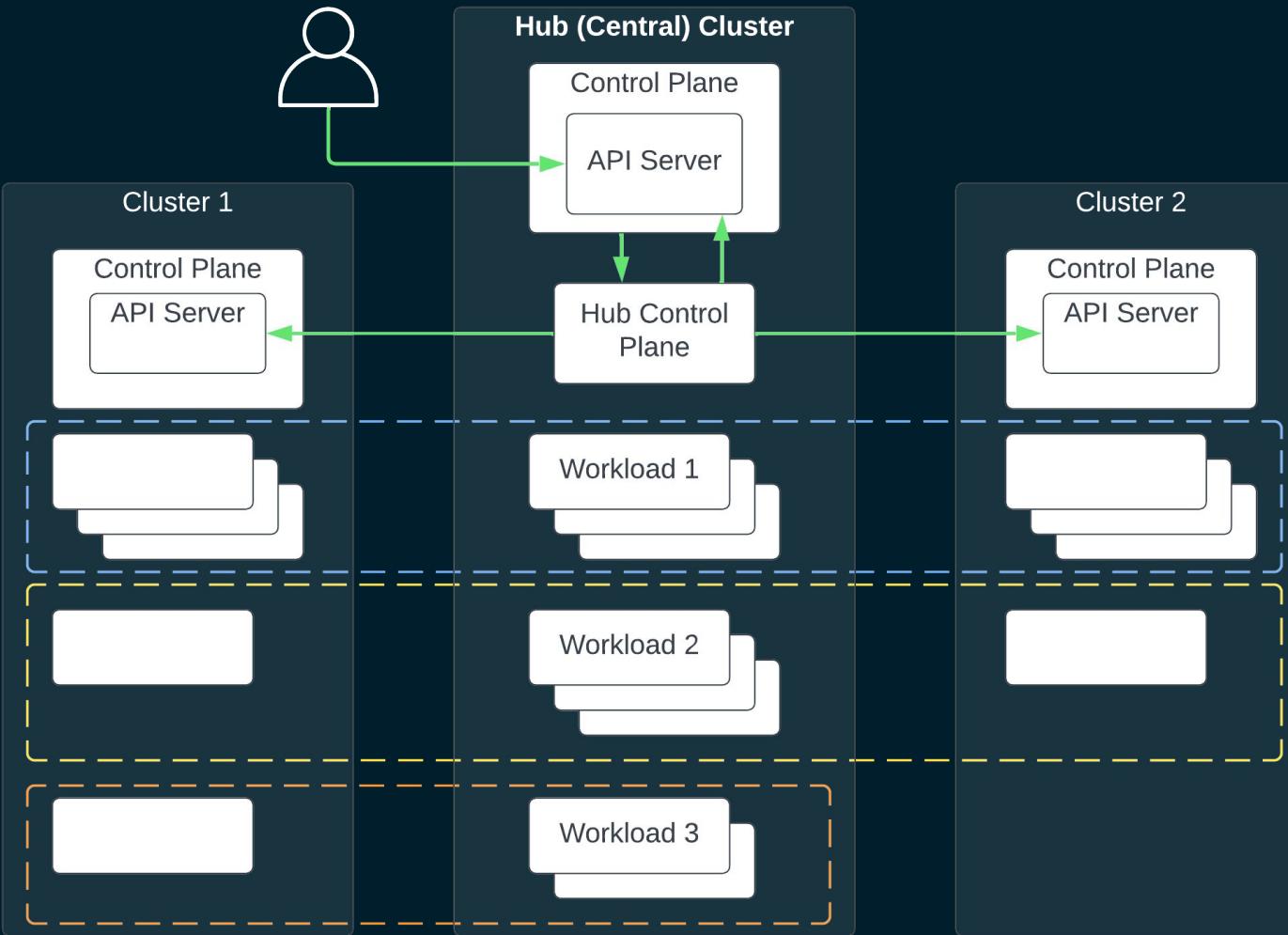
MongoDB Cluster Inventory

- Creation of RBAC on member clusters
- Store the kubeconfig of member clusters

```
kubeconfig:  
  apiVersion: v1  
  clusters:  
    - cluster:  
        certificate-authority-data: <cluster-1-ca.crt>  
        server: https://127.0.0.1:41727  
        name: cluster-1  
  contexts:  
    - context:  
        cluster: cluster-1  
        namespace: mongodb  
        user: cluster-user-1  
        name: cluster-1  
  kind: Config  
  users:  
    - name: cluster-user-1  
      user:  
        token: <cluster-1-token>
```

Workload Distribution





Considerations:

- Redundancy
- Geographical
- Resources

Workload Distribution

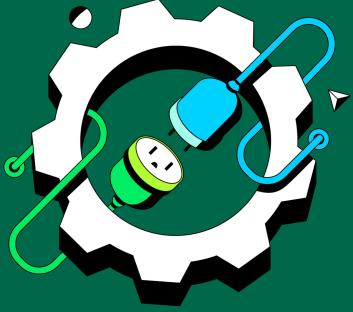
- Declarative approach using *spec.clusterSpecList*
- Users request distribution of the database nodes
- In case of Kubernetes cluster failure, we can re-distribute statefulsets to healthy clusters

```
1 apiVersion: mongodb.com/v1
2 kind: MongoDBMultiCluster
3 metadata:
4   name: multi-replica-set
5 spec:
6   version: 6.0.6-enterprise
7   type: ReplicaSet
8   persistent: false
9   duplicateServiceObjects: false
10  credentials: my-credentials
11  opsManager:
12    configMapRef:
13      name: my-project
14  clusterSpecList:
15    - clusterName: cluster1.mongodb.com
16      members: 2
17      statefulSet:
18        spec:
19          # to override the default storage class for the pv
20          volumeClaimTemplates:
21            - metadata:
22              name: data
23            spec:
24              accessModes: [ "ReadWriteOnce" ]
25              storageClassName: "gp2"
26    - clusterName: cluster2.mongodb.com
27      members: 1
28      statefulSet:
29        spec:
30          volumeClaimTemplates:
31            - metadata:
32              name: data
33            spec:
34              accessModes: [ "ReadWriteOnce" ]
35              storageClassName: "gp2"
36    - clusterName: cluster3.mongodb.com
37      members: 1
```

Networking



Networking Considerations



Network topology



Security



Performance

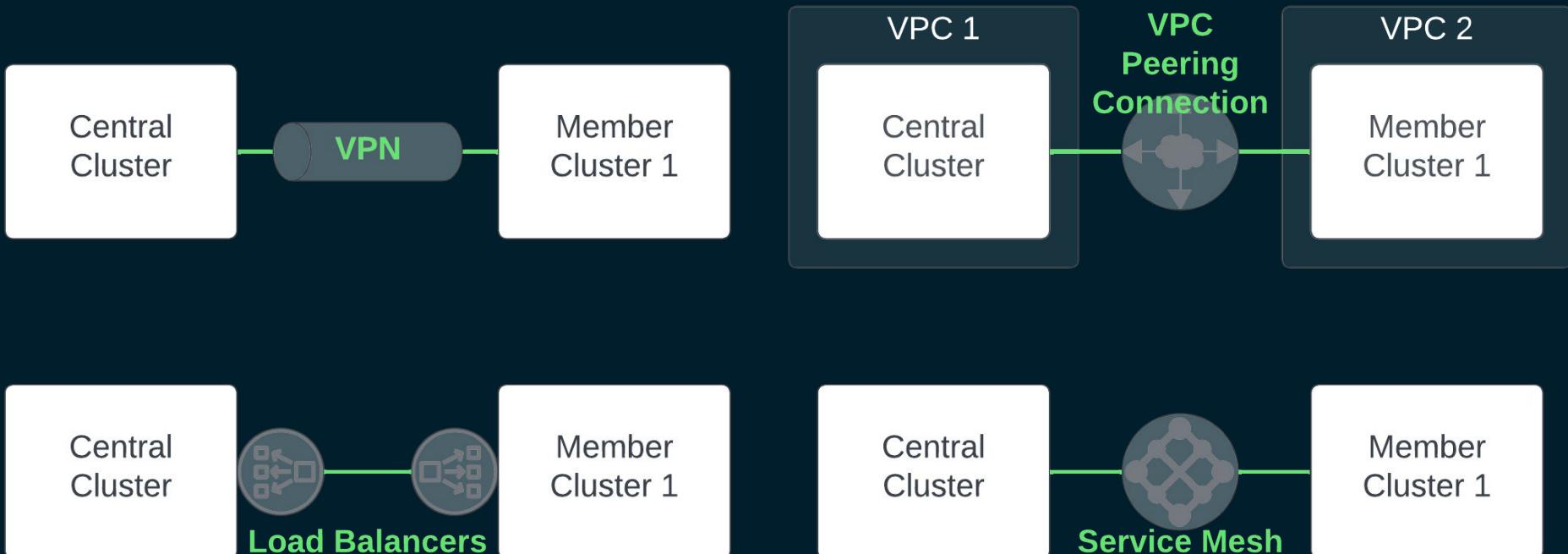


Reliability & resilience



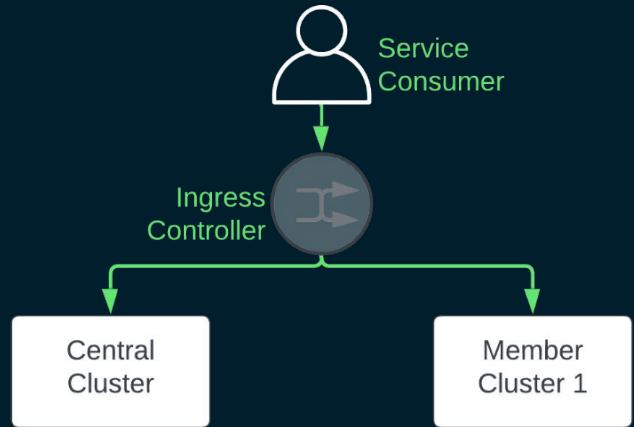
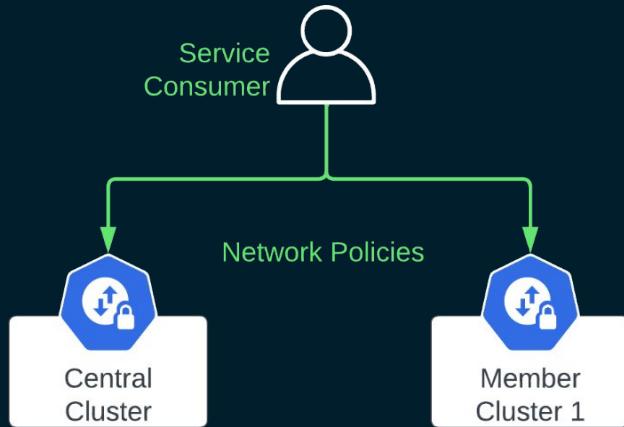
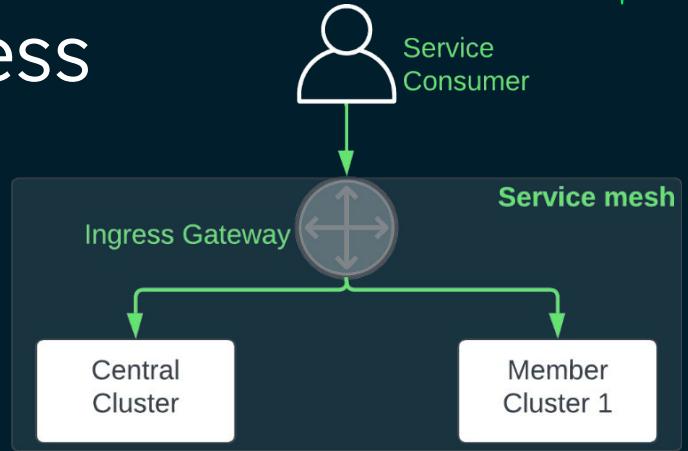
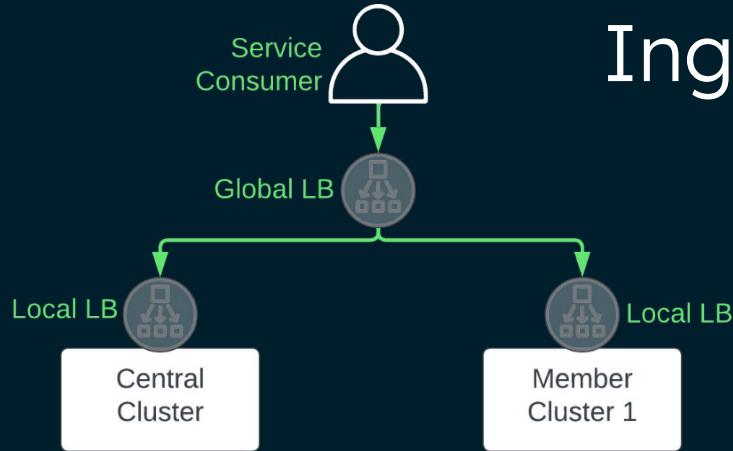


Inter-Cluster Networking

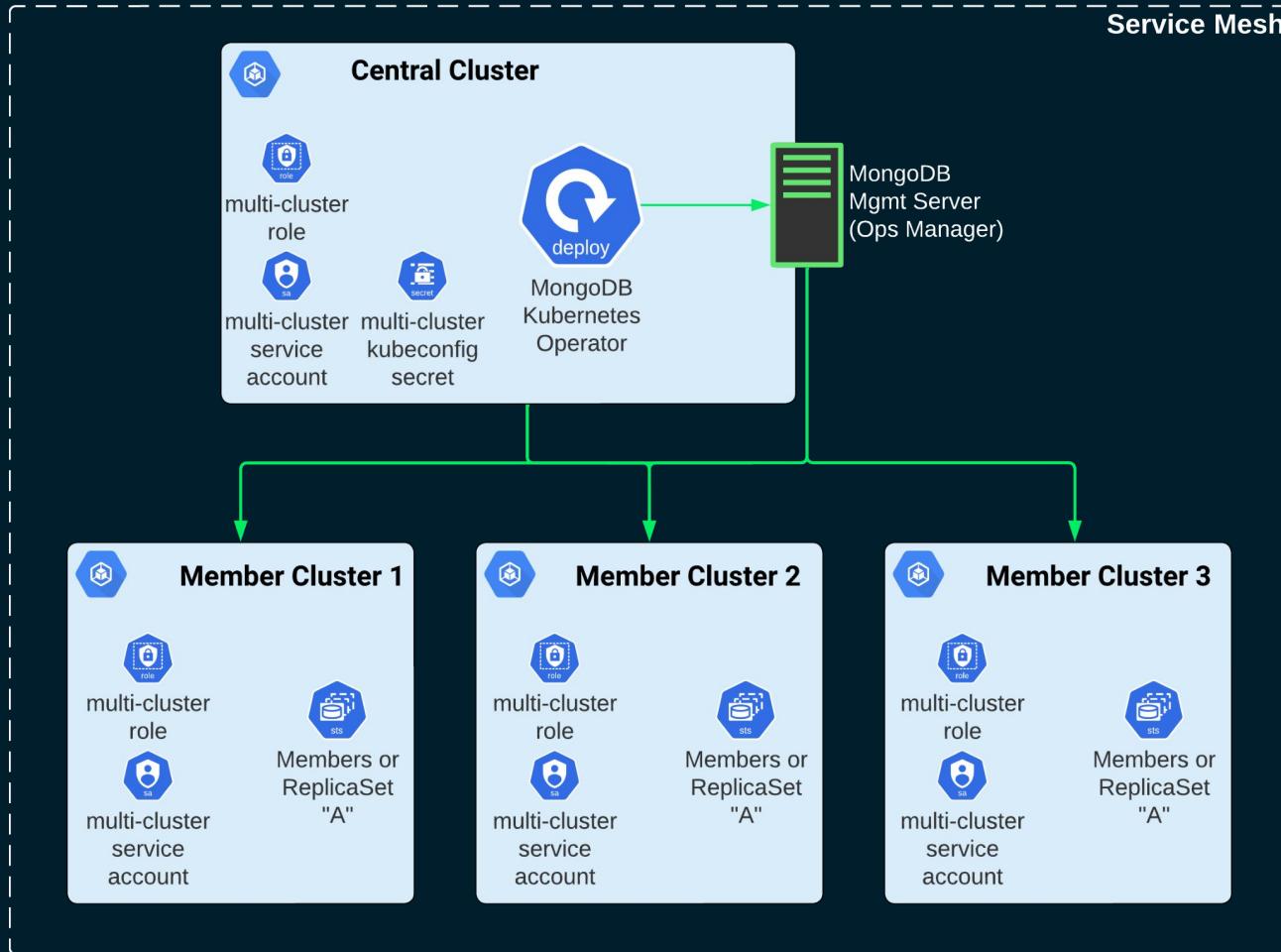




Ingress & Egress



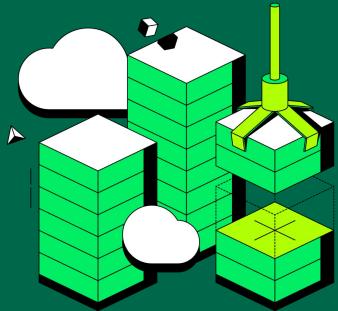
Using a service mesh



Multi-Cluster Controllers and Operators



Controller & Operator Design Considerations



Consistent state
across clusters



Resiliency -
specifically to cluster
failure



Scalability



Security





Functional Needs

- Cluster Inventory & Workload Distribution
- Failure handling
- Monitoring and Observability



Existing Solutions





Karmada

Considers cluster availability and actual resources

Automated or manual cross-cluster failover

Cross-cluster service discovery

Uses Submariner to establish inter-cluster connectivity

More declarative

Optional karmada-agent on each member/spoke clusters

Can retrofit a single cluster workload to multi-cluster



Open Cluster Management

Considers cluster availability and actual resources

Automated or manual cross-cluster failover

Cross-cluster service discovery

Uses Cluster Proxy to establish L4 connectivity, or DNS integration

More prescriptive

Klusterlet agent on each member/spoke clusters



MongoDB's Journey

- Workload distribution: Prescriptive → Declarative
- Resource Awareness
- Automated failover → More control
- Operator High Availability



Future of Multi-Cluster





Future of Multi-Cluster

- More powerful & user-friendly multi-cluster management tools
- Increased use of AI and ML to automate tasks
- Better network support for multi-cluster Kubernetes
- Better support for hybrid and multi-cloud deployments



Q&A

Dan Mckean

Senior Product Manager, Kubernetes

George Peter Hantzaras

Director of Engineering, Kubernetes



[linktr.ee/kubecon_multi
cluster_patterns](https://linktr.ee/kubecon_multi_cluster_patterns)