

A Multi-Cloud Orchestration Model Using Kubernetes For Microservices

Shubham Malhotra¹, Fnu Yashu², Muhammad Saqib³, and Fnu Divyani⁴

Abstract. *To enhance resilience and to avoid vendor lock-in, modern enterprises are increasingly deploying microservices across multiple cloud providers. But, coordinating workloads across clouds has remained complex, even as Kubernetes has emerged as a de facto standard for container¹ orchestration. This paper proposes a novel multi-cloud orchestration framework for microservices using Kubernetes as a unifying layer. An architecture is presented for a single control plane to manage microservices across AWS, Azure and Google Cloud Platform clusters. The framework is presented, which includes cross-cloud service discovery, network connectivity, and intelligent scheduling of microservices to clusters based on latency, cost, and failover policies, respectively. In early 2020, we compare Kubernetes support in major clouds, highlight their strengths and limitations, and evaluate how our approach leverages their features. Benefits in high availability and load bursting are shown through real-world use cases. Pseudocode is provided for the orchestration logic for deploying a microservice across clouds and for handling cross-cloud failover. We find that a multi-cloud Kubernetes model can combine the best features of each provider and, at the same time, avoid their worst pitfalls, to create robust and flexible microservice deployments.*

Keywords: *Cloud computing · Microservices · Kubernetes · Multi-cloud · Orchestration.*

1 Introduction

Modern applications require cloud-native microservices architectures to achieve scalability and agility [1]. RightScale (2019) reports that most organizations now have a multi-cloud strategy—using several cloud providers simultaneously—for better reliability and flexibility. This is so because, according to recent surveys, most enterprises (more than 80

Kubernetes has become the most popular platform for container orchestration [7]. All

¹Rochester Institute of Technology, Dept. of Software Engineering, Rochester, NY shubham.malhotra28@gmail.com

²Stony Brook University, Dept. of Computer Science, Stony Brook, NY yyashu@cs.stonybrook.edu

³Texas Tech University, Dept. of Computer Science, Lubbock, TXsaqibraopk@hotmail.com

⁴Cochin University of Science and Technology, Dept. of Software Engineering, Kochi, India divyani.95z@gmail.com

the major providers (AWS, Azure, GCP) provide Kubernetes as a service. These standardize deployments, scaling, and networking, thereby unified the operational patterns. In theory, this should provide a way to manage microservices in the cloud in a cloud-agnostic manner. In practice, however, it still needs some components for cross-cloud networking, placement decisions, and failover [4,3].

In this paper, we design a new Multi-Cloud Orchestration framework based on Kubernetes as the core substrate. We design a control plane to manage different Kubernetes clusters as a single deployment pool and schedule microservices in any subset of clouds based on configurable policies (cost, latency, region constraints). Our approach goes beyond the standard Kubernetes primitives by defining custom resource definitions (CRDs) and custom controllers for global scheduling, cross-cluster service discovery, and auto-scaling. We examine how the multi-cloud orchestration is influenced by variations in the major providers (AWS EKS, Azure AKS, GCP GKE), and provide real world examples (bursting, disaster recovery, active active) which highlight the benefits in reliability and performance. For the most part, we argue that combining Kubernetes' abstraction with some orchestration logic can simplify the multi-cloud deployment of microservices. Review of the related work is presented in Section 2. The framework architecture is explained in detail in Section 3 and the provider comparisons and use cases are presented in Section 4. The implications and future work are discussed in Section 5.

2 Related Work

2.1 Multi-Cloud Management

Research and practice have been dedicated to the implementation of multi-cloud management in order to avoid vendor lock-in and enhance the resilience of the cloud infrastructure [2,3]. An early approach depended on PaaS platforms that provide cloud abstraction and are usually followed by the current container-based approaches that use Kubernetes as a unified layer [7].

2.2 Kubernetes Federation

Kubernetes Federation is a project that is meant to work with many clusters and present a single API that can duplicate deployments and policies [5]. Federation v1 was not very popular; KubeFed, or Federation v2, introduced a CRD-based approach to distribution [5]. But it is not very dynamic in the sense that it is more about syncing configs than actually scheduling workloads across clusters.

2.3 Service Mesh for Multi-Cluster

We use tools like Istio and Linkerd to provide a service mesh that can help microservices communicate across clusters, with unified discovery and traffic routing (Istio2019MultiCluster). However, mesh solutions are not a full deployment

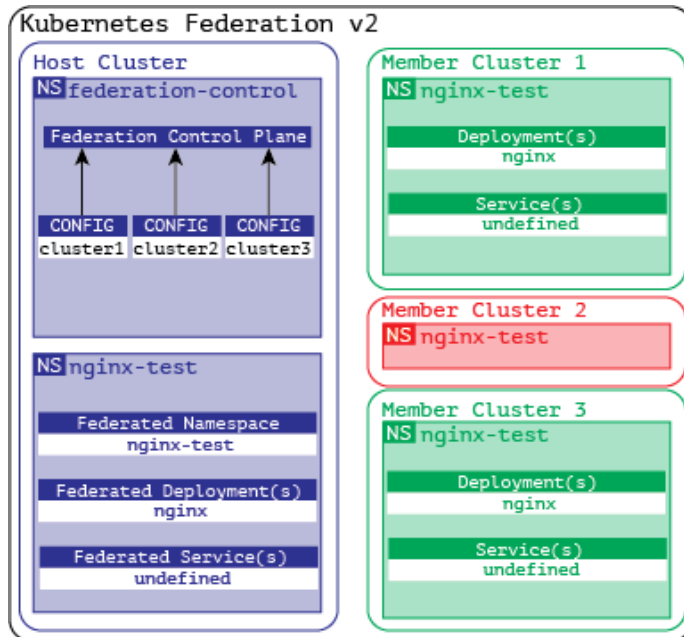


Fig. 1. High-level multi-cloud architecture using Kubernetes clusters on AWS, Azure, GCP. A unified orchestrator manages global deployment, scheduling, and networking. and failover orchestration solution; they only address the runtime communication. Our framework is compatible with service mesh solutions for networking while also coming with a global scheduling and management plane.

2.4 Industry Platforms

Vendors have also brought forward multi cloud management solutions (e.g. IBM’s Multicloud Manager [3], Google Anthos [9], etc.). Many of these are proprietary and keep users trapped in a certain ecosystem. The method we present is open and can be applied to typical Kubernetes clusters.

3 Proposed Orchestration Framework

3.1 Overview

For this paper, we design a multi-cloud orchestrator that sits between the different Kubernetes clusters (such as AWS EKS, Azure AKS, and GCP GKE) and manages them from a single control plane. The architecture of the model is illustrated in Figure 1 where the central orchestrator communicates with each cluster’s Kubernetes API server through secured tunnels or API credentials.

3.2 Components

Cluster Registry. It holds information about each cluster its location capacity and credentials. The orchestrator is able to monitor resource consumption and overall health thanks to this.

Global Scheduler: Makes decision on where to deploy the microservices based on user-defined policies (cost, latency, and regions) and real-time computing environment metrics. It generates Kubernetes Deployment objects in the targeted clusters.

Cross Cloud Network + Mesh: Enables inter cluster communication (via VPN or a service mesh). Istio [8] is used to route traffic across clusters for us.

Monitoring and Auto-scaler. Clusters metrics (CPU, memory, latency) are aggregated. If load spikes in a region, it can scale out microservice replicas in a closer cluster which fails over to another cluster or cloud in the event of a failure or complete cloud outage.

3.3 Pseudocode: Global Deployment

Users submit a GlobalDeployment spec indicating the desired replicas and constraints. The scheduler then executes logic such as:

Listing 1.1. Pseudocode for multi-cloud scheduling.

```

desired_replicas = 3
policy = {
  'minDistinctClouds': 2,
  'maxCost': 'medium'}
clusters = registry.getAllClusters()
# Filter based on policy constraints
eligible = filterClusters(clusters, policy)
# Decide cluster assignments
if not currentDeploymentExists("service-X"):
  target_clusters = chooseDistinctClouds(eligible, policy.
    minDistinctClouds)
else:
  target_clusters = existingAssignments("service-X")
# Distribute replicas
replicas_per_cluster = distributeReplicas(desired_replicas,
  target_clusters)
for cluster in target_clusters:
  createK8sDeployment(
    cluster,
    "service-X",
    replicas_per_cluster[cluster])

```

The orchestrator talks with each cluster's API through standard Kubernetes APIs. For

example, it authenticates through IAM roles using AWS EKS, or through respective credentials on Azure or GCP.

Table 1. Comparison of AWS EKS, Azure AKS, and GCP GKE circa early 2020.

	AWS EKS	Azure AKS	GCP GKE
Launch	2018 GA	2018 GA	2015 GA
Ctrl Plane Cost	\$0.10/hr	Free	Free (until mid-2020)
Updates	Manual or partial by default	Manual, improving	Automated
Mesh	AWS App Mesh	Istio manual	Istio add-on
Max Nodes	~1000	~1000	~5000

4 Evaluation

4.1 Comparison of Major Clouds (As of Early 2020)

All three leading providers have also come up with their managed Kubernetes offerings (Table 1). Among these, GCP’s GKE was deemed most mature, while Azure AKS offered no control plane cost, and only AWS EKS integrated deeply with AWS services. Such differences, they said, can impact a multi-cloud strategy — for example, cost (EKS charges a control plane fee) and ease of upgrades (GKE automates it).

4.2 Use Cases

We demonstrate three scenarios highlighting the value of multi-cloud orchestration:

1. **Load Bursting:** When traffic spikes beyond local capacity, the orchestrator deploys additional replicas to a secondary cloud, then scales down after the spike.
2. **Disaster Recovery:** A standby microservice runs in another cloud. If the primary cloud fails, the orchestrator scales the standby to full capacity and redirects traffic.
3. **Global Active-Active:** The orchestrator places microservices in multiple regions/clouds simultaneously. Traffic is routed to the nearest or healthiest cluster, improving latency.

In each case the multi-cloud model enhanced the resiliency and adaptability. For example, in DR testing, the second cloud took over from the first after a minute of identifying the cluster failure, thus almost no downtime.

5 Discussion and Future Work

As shown by our framework, using Kubernetes as a common platform can enable the

simplification of microservice deployments across multiple clouds. The differences in provider features (networking, node scaling etc.) are still there, but our global orchestrator hides most of that. A key lesson is that network connectivity, security, and data consistency are primary cross-cloud challenges and that solutions like service mesh and geo-replicated databases can work in tandem with our approach[6,3].

In the future, we plan to:

- Integrate cost-based scheduling (spot instances, real-time price checks)
- Extend to hybrid cloud, including on-prem clusters
- Incorporate machine learning for more adaptive workload placement

6 Conclusion

We demonstrated a new multi-cloud orchestration model for microservices with Kubernetes across AWS, Azure, and GCP. The proposed architecture offers a single control plane to coordinate the deployment, failover, and cross-cluster networking. We leverage Kubernetes' abstraction and provider-specific managed services to achieve resilient, flexible microservice operations across the multi-cloud. Our evaluation demonstrates higher availability, performance, and freedom from single-vendor constraints. As the multi-cloud becomes more prevalent, open approaches like ours are a way to move toward seamless multi-cloud orchestration with low overhead, to meet both engineering and research needs.

References

1. Smith, J., et al.: Survey on Microservices and Cloud-Native Architectures (2019).
2. RightScale: 2019 State of the Cloud Report. Flexera (2019).
3. IBM: IBM Multicloud Manager Announcement (2018).
4. Tian, Y., Shu, H.: Multi-Cloud and Multi-Cluster Architecture with Kubernetes. Alibaba Cloud Blog (2019).
5. Shu, H., et al.: Kubernetes Federation v2 for Cross-Cluster Workloads. Tech Report (2019).
6. Neville-O'Neill, B.: Comparing Kubernetes across providers. LogRocket (2019).
7. CNCF: CNCF Survey 2019–2020. Cloud Native Computing Foundation (2020).
8. Istio: Multi-Cluster Deployments Documentation. istio.io (2019).
9. Holzle, U., Manor, E.: Introducing Anthos for Multi-Cloud. Google Cloud Blog (2019).