



---

## Kubernetes for Multi-Cloud and Hybrid Cloud: Orchestration, Scaling, and Security Challenges

**Bhanuprakash Madupati**

MNIT, MN

---

**Abstract:** Kubernetes enables organizations to manage applications globally, from multi-cloud to hybrid clouds. This paper presents challenges and solutions for workload orchestration, scaling, and security in such infrastructures. Running workloads across different cloud providers enables one to escape vendor lock-in and improve service availability. However, it leads to more complexity in terms of management and security. This paper will analyze the tendency of Kubernetes to orchestrate across multiple clouds in which it operates in a Multi-Cloud environment. It covers important solutions for scaling & managing cross-cloud workloads and securing applications in a multi-cloud or hybrid-cloud environment. Security discussions in the paper include managing identity and access across multiple cloud providers and supporting zero-trust security models. We then examine optimal Kubernetes operations across multi-cloud and hybrid cloud infrastructures, considering solutions like Cluster Federation, service mesh and encryption methods.

**Keywords:** Kubernetes, Multi-Clouds, Hybrid Cloud, Workloads Scheduling/Scaling more secure.

---

### 1. Introduction

Cloud computing is one of the fundamental building blocks for modern IT infrastructure. It enables organizations to use scalable resources, reduce costs, and grow big or shrink easily. In a world where businesses more commonly leverage multiple cloud providers (multi-cloud) and private on-premise systems integrated with public clouds (hybrid cloud), Kubernetes has become an increasingly important tool for enterprises to successfully manage such complex environments. Kubernetes abstracts the underlying infrastructure, provides primitives for defining microservices and their deployment patterns, and automates the management of these microservices across various clouds [1].

Multi-cloud and hybrid cloud architectures provide benefits such as overcoming vendor lock-in (lock of a customer to the products or services of one provider) and adding redundancy [3]. But it also comes with challenges—including the need for workload orchestration, serviceability and security across multiple platforms. It solves all these problems by providing cluster federation capabilities, automatic scaling and security policies, which makes it a most suitable choice for managing containerized applications across clouds [5]—areas of focus — Workload orchestration and scaling in multi-cloud and hybrid cloud; Challenges related to security.

Section 2 describes how Kubernetes provides tools and methods for orchestrating workloads across clouds between them [4]. Section 3 discusses difficulties and approaches when scaling applications in hybrid & multi-cloud environments using Kubernetes [5]. Lastly, Section 4 covers security challenges unique to such environments and highlights guidelines for load securing in these contexts [2], [7].

Kubernetes provides an orchestration Layer to manage your workloads across different cloud providers and environments, which is important for managing modern Cloud Infrastructure. Here, we are going to experiment with these topics in detail and try to understand what Kubernetes can do and where it is a blocker in a multi-cloud and hybrid-cloud environment.

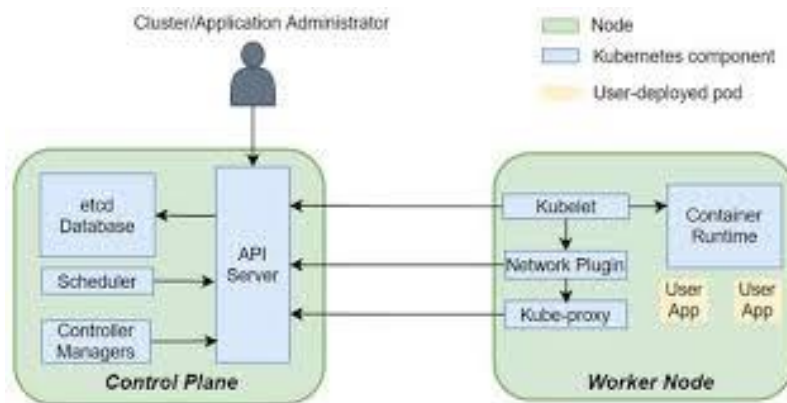


## 2. Multi-Cloud/Hybrid Cloud Workload Orchestration

Kubernetes manages workloads deployed across multi-cloud and hybrid-cloud environments. Its main benefit is that it abstracts infrastructure concerns and allows users to move workloads seamlessly across different cloud platforms. Kubernetes enables multi-cloud or hybrid cloud architecture with its key components and tools in workload management, data synchronization, and networking.

### Cluster Federation

One of the important features in Kubernetes is called Cluster Federation (also known as Kubernetes Federation), which solves the problem by managing multiple Kubernetes clusters across different cloud environments as a single entity. Federation allows scheduling workloads across clusters and distributing them among cluster members, augmenting resource utilization while making applications resilient to single-cluster failures. Federation enables an organization to deploy workloads in different geographic regions, which results in greater availability and lower latency [3].



**Figure 1:** Diagram showing Kubernetes cluster federation across multiple clouds, with a visualization of how workloads can be distributed across different regions (public clouds and private clouds)

Not only can you drastically simplify the distribution of your workloads across Kubernetes clusters with Kubernetes Federation, but it also gives you an effectively free disaster recovery/failover mechanism in case a cluster goes down as well — if a cluster happens to go away in cloud A, the workload is automatically moved over to another live cluster in cloud provider B. Highlighted in hybrid cloud settings, as on-site resources are glued to the public services in the cloud that along with ensuring continuous operations [1].

### Multi-Cloud Networking

However, maintaining communication between clusters is a major problem in managing workloads across multiple clouds. To solve this issue, Kubernetes uses both Istio and Calico for service discovery, load balancing, and traffic routing between clouds. Istio is a service mesh architecture that makes it easier to manage microservices [4] and provides a uniform way to secure, connect, and monitor microservices across multiple clusters.

**Table 1:** Comparison of networking solutions in Kubernetes multi-cloud orchestration.

Networking Tool	Use Case	Key Features
Istio	Service Mesh for Microservices	Traffic routing, security, observability
Calico	Network Security	Fine-grained control, secure communication
Flannel	Simple Networking for Kubernetes Clusters	Flat networking for inter-pod communication

### Data Synchronization Across Clouds

Data consistency in multi-cloud environments is challenging, given the resources' distributed nature. Rook and Ceph - Kubernetes provides the mechanism for managing data across clouds via tools like Rook and Ceph. In this regard, such a cloud-native storage solution allows provision and management of persistent storage across different clouds and workloads using the Stateful Set to access consistent data regardless of the cluster location [1], [7].

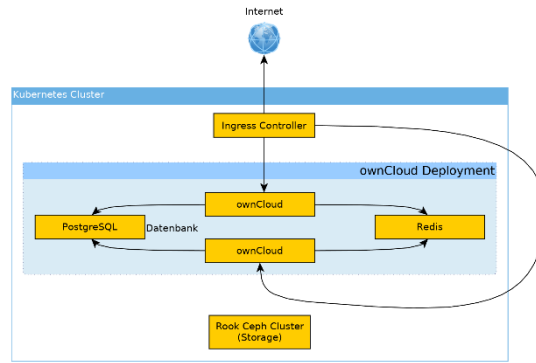


Figure 2: Diagram showing data synchronization using Rook/Ceph across multiple Kubernetes clusters in different cloud environments.

As Stateful applications are becoming popular in Kubernetes, it is crucial now that the data in one cloud environment where the application has been deployed is the same as in another cloud environment. Rook For example, with Rook, it is possible to set up distributed storage systems to guarantee that data replication and consistency are kept and to ensure the absence of data non-availability when a cluster fails [7].

**Unified Management Through Declarative Configurations**

With Kubernetes, we use a declarative configuration (a configuration that specifies what you want to exist and not how it should be done) of our desired state for Kubernetes to manage our deployments. This approach makes it easier to manage multi-cloud and hybrid cloud environments as administrators can define the infrastructure and workload configuration once and deploy them all the same across clouds [5].

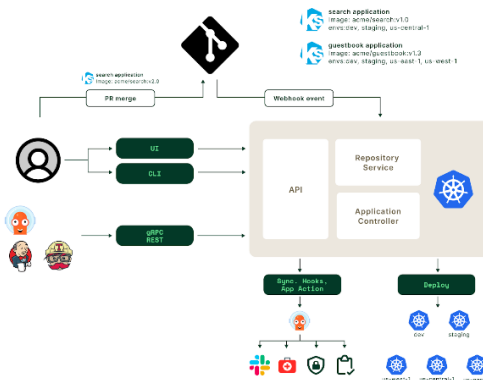


Figure 3: Diagram of Kubernetes' declarative configuration process, showing how desired states are applied across multiple clusters in multi-cloud environments.

The declarative method facilitates workload management across complicated cloud installations by minimizing human error and guaranteeing consistent and reliable infrastructure.

**3. Scaling Kubernetes in Multi-Cloud/Hybrid Cloud**

Scaling is one of the most important features for managing multi-cloud or hybrid-cloud applications. Kubernetes scaling features are part of the lifecycle management functionality, providing a powerful scale in/out capability and automatically adjusting workloads dynamically based on the resource demand, one of the key areas Kubernetes helps with when working in highly dynamic distributed environments like multi-cloud. In this post, we will discuss how Kubernetes manages Scaling, the horizontal and vertical configuration of the pods, and problems regarding load balancing across clouds. We will also discuss how to minimize latencies in the case of hybrid architectures.

**Horizontal & Vertical Scalability**

Kubernetes supports horizontal Scaling, which adds more pods when the load increases, and vertical Scaling, which assigns the CPU or memory resources to existing pods. This includes using Kubernetes' autoscaling

features to manage workloads that stretch across multiple cloud providers and in multi-cloud and hybrid cloud environments.

Horizontal Scaling is mostly used with the help of Horizontal Pod Autoscaler (HPA). The controller manages the deployment, automatically scaling the number of pods based on average observed CPU utilization (or anything else if I create custom metrics to expose). HPA can even span multiple clusters in multi-cloud configurations, enabling a single scaling logic to be applied across clouds [5].

**Vertical Scaling:** Kubernetes can change the amount of resources allocated to our existing containers by increasing or decreasing the CPU or memory because more components are needed due to workload needs. This is orchestrated by the Vertical Pod Autoscaler (VPA), which can be handy in places where increasing resources scales much higher than adding another pod [2], [5].

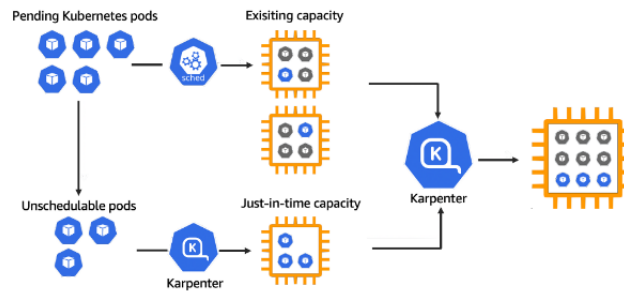


Figure 4: Diagram illustrating Horizontal and Vertical Scaling in a Kubernetes environment, showing how workloads adjust across multi-cloud clusters.

**Balancing loads across clouds**

In multi-cloud architecture, the main reason for distributing traffic across clusters is that we want workloads to be balanced and not to leverage too much workload on one particular resource. Through Ingress Controllers and depending on the actual cloud provider, Kubernetes can provide both intra-cloud and inter-cloud load balancing capabilities [5].

Kubernetes is always on, with Global Load Balancing, which distributes traffic to the closest cluster based on global proximity or latency. In such an environment, we often use tools like NGINX Ingress Controller or Cloud-Native solutions like Google Cloud Load Balancer to efficiently divide traffic.

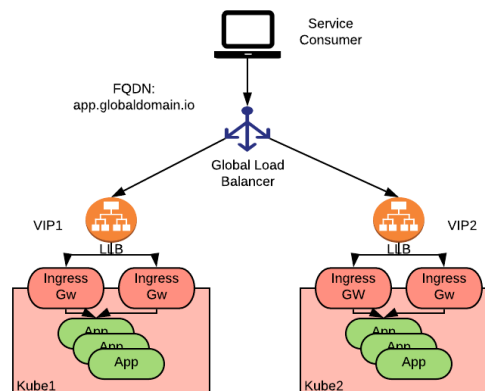


Figure 5: Diagram of Global Load Balancing using Kubernetes, demonstrating traffic routing between geographically distributed clusters in multi-cloud environments.

**Managing Latency in Multi-Cloud and Hybrid Cloud**

They ensure low latency to all components, especially in hybrid cloud architectures where workloads are distributed between on-premise infrastructure and public clouds. For example, Kubernetes helps you reduce latency by sub-milliseconds to milliseconds by scheduling the workloads close to data sources or a specific

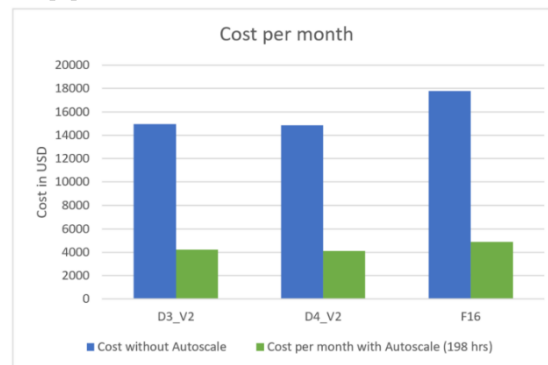
region serving high demand. This technique, which is also called edge computing, guarantees that the end-user experience will remain the same even when workloads are extended through multiple regions or clouds [3], [4]. This network proximity control ensures that workloads (e.g., video streaming, financial services) requiring real-time processing are placed in ultra-low latency regions.

**Table 2:** Latency Management Techniques in Multi-Cloud/Hybrid Cloud Environments.

Technique	Description	Use Case
<b>Edge Computing</b>	Workloads placed near the data source or user	Low-latency applications (e.g., streaming)
<b>Global Load Balancing</b>	Routes traffic based on geographical proximity	Applications with users in multiple regions
<b>Multi-region Deployment</b>	Deploy workloads across multiple regions to reduce lag	High-availability workloads

### Autoscaling and Cost Optimization

Autoscaling Resources: One of the core advantages and original target use cases of Kubernetes, within multi-cloud/hybrid-cloud environments, is that you can configure your resources to scale up when demand increases to reduce costs. Depending on the needs, the Cluster Autoscaler can resize Kubernetes Clusters dynamically by adding or deleting nodes seamlessly across clouds to manage workloads at a lower price. This is especially useful in hybrid cloud environments where resources can be scaled down locally off-peak charge and public cloud resources for peak requests [7].



**Figure 6:** Graph showing the cost benefits of autoscaling across hybrid cloud environments, demonstrating how scaling up and down based on demand can lead to cost savings.

### 4. Kubernetes Security Challenges in Multi-Cloud / Hybrid Cloud

Running security in the multi-cloud and hybrid cloud settings is even more difficult – because each platform supports separate security models, has distributed resources and —just like with email— all your relevant processes are desperate for some consistent identity management. Kubernetes provides a plethora of tools and mechanisms to tackle these challenges, ranging from Role-Based Access Control (RBAC) for authorization to service mesh architectures [1][2] like Istio for securing service-to-service communications. This part covers the primary issues and better practices in Kubernetes for multi-cloud/hybrid cloud deployments.

#### Identity and Access Management (IAM) in Multi-Cloud

Most multi-cloud providers will use separate IAM (Identity and Access Management) frameworks for their cloud environment running in the data centre. Integrating frameworks and ensuring the security level is maintained across all platforms can get some work. One of the ways that Kubernetes secures this access is through Role-Based Access Control (RBAC) to control who can do what within a cluster, giving fine-grained controls over who can see and change which resources [2].

RBAC provides the toolbox for administrators to define roles that constrain what users can do within specific namespaces in a Kubernetes cluster, whether the cluster is deployed on a public cloud or on-premise in a hybrid cloud setting. Having these roles consistently applied across clusters on different cloud providers is still very important to prevent unauthorized access [3].



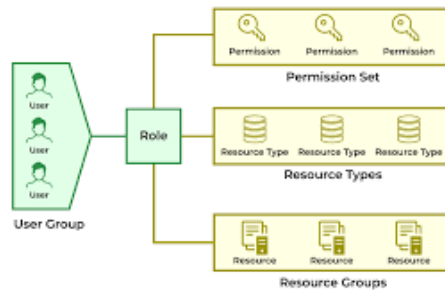


Figure 7: Diagram of Kubernetes RBAC, showing role assignments across multiple clusters in a hybrid cloud environment.

Service accounts and Kubernetes Namespaces can also be used to isolate workloads in a single cluster, limiting services and users to what resources they have access to.

**Zero Trust Security Models**

The zero-trust security model is based on decentralized trust: it expects all requests—outside or inside—to be checked before access is permitted. This model is directly enabled by Kubernetes and tools like Istio and Linkerd, which offer mutual Transport Layer Security (mTLS) to secure inter-service communications with strong identity validations [4], [5].

In multi-cloud scenarios with workloads distributed across multiple regions and cloud providers, these policies can be enforced uniformly because Istio secures communication between services regardless of geographical distribution. This means that even if the security of one cloud provider is busted, others should be safe as long as strong authentication and encryption practices prevent hackers from accessing user data in plain text.

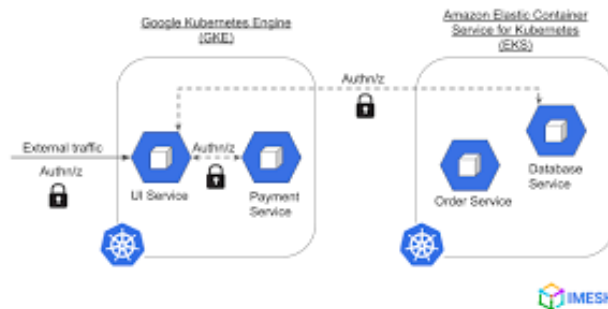


Figure 8: Diagram illustrating the Zero-Trust model in a Kubernetes multi-cloud setup, showing how mTLS encrypts communication between microservices across clouds.

**Network Security & Policies**

Another equally important challenge is guaranteeing secure network communication across multi-clouds. Kubernetes network policies: While Kubernetes makes sure that only intended traffic can arrive at the pod, it also provides network policies for administrators to control what is allowed between pods, between services, and with external entities. Kubernetes applies service-to-service isolation via well-defined NetworkPolicy resources that describe which services are allowed to talk to which [5].

Other tools like Calico and Cilium extend these inbuilt network security capabilities. They offer various advanced features, such as data-in-motion encryption and multi-cloud intrusion detection against Kubernetes. These tools are necessary for secure communication across geo-distributed clusters in a multi-cloud.

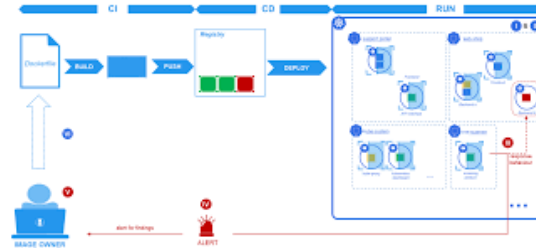
**Table 3:** Comparison of Kubernetes Network Security Tools for Multi-Cloud.

Tool	Feature	Use Case
Calico	Fine-grained control, encryption	Securing communication between clusters
Cilium	Data encryption, monitoring	Intrusion detection, secure communication
Network Policy	Pod-level traffic isolation	Defining internal and external traffic rules

### Data Encryption and Vulnerability Management

Encryption at rest and in transit: Data security for multi-cloud and hybrid cloud environments can not be compromised. Kubernetes supports many cloud-native encryption services, such as AWS KMS, Google Cloud KMS, and Azure Key Vault [2]. These services are offered to ensure that even if data is distributed over multiple cloud providers, it will still be encrypted, and this sensitive information will remain secure.

Security: Besides encryption, you must also manage the vulnerabilities in container images. Services like Aqua Security and Anchore can scan container images for vulnerabilities before they are introduced to Kubernetes clusters. This helps to solve any types of would-be security issues before they can get into production [7].



**Figure 9:** Flowchart illustrating container vulnerability management in Kubernetes, showing how images are scanned before deployment and securely stored.

### Compliance & Regulatory difficulties

Multi-cloud environments are not only cumbersome to manage because GDPR, HIPAA, and PCI-DSS can make heading your way into data laws for different locations time-consuming based on the geographic location of the data. With Kubernetes, compliance rules can be enforced by restricting workloads to stay within certain specified regions/cloud providers based on what is legally required due to local regulations [7].

Node selectors and taints enable Kubernetes administrators to prevent data processing in regions that do not follow local laws by deploying sensitive workloads only in specific clouds or regions.



**Figure 10:** Diagram showing Kubernetes workload placement based on regulatory requirements, with node selectors ensuring workloads are restricted to specific geographic regions.

### 5. Conclusion

1. This is accomplished by Cluster Federation, which helps in cross-cloud management of clusters, where the same application can be deployed seamlessly on multiple cloud and/or on-prem clusters.
2. Kubernetes's horizontal and vertical scaling mechanisms make it easy to scale, enabling efficient resource consumption and saving costs when scaled up. Global load balancing streamlines traffic management to cut down on latency and improve the availability of web services.
3. Improved security for Multi-Cloud Kubernetes Environments: Enabled with resources RBAC, Zero-trust security model and Network Security Policies. Tools like Istio and Calico for service communication and traffic management provide advanced security.
4. Kubernetes and its storage ecosystem, represented by Rook and backed by Ceph, are essential in solving the data synchronization and storage management problem when achieving consistent availability across cloud environments, thus serving stateful workloads.



5. Data privacy and regional regulation compliance: This capability prevents your workloads from being deployed in a given geographic location in the case of data leaks or when such movements are not allowed for reasons like GDPR.
6. Multi-cloud Kubernetes management best practices include using declarative configurations, integrating vulnerability scanning tools, and enabling cloud-native encryption solutions for data at rest/in transit.

## Reference

- [1]. P. Buchner, "From the Cloud to the Edge: An Infrastructure for Cloud & Edge Computing," Epub.jku.at, 2019. [Online]. Available: <https://epub.jku.at/obvulihs/content/titleinfo/4382124>.
- [2]. T. Autio, "Securing a Kubernetes Cluster on Google Cloud Platform," Urn.fi, 2021. doi: <http://www.theseus.fi/handle/10024/501918>. Available: <https://urn.fi/URN:NBN:fi:amk-2021060314068>.
- [3]. H. Mfula, A. Ylä-Jääski, and J. Nurminen, "Seamless Kubernetes Cluster Management in Multi-Cloud and Edge 5G Applications," Aug. 2021. [Online]. Available: [https://tuhat.helsinki.fi/ws/portalfiles/portal/167373071/\\_harrison\\_mfula\\_ccaaS\\_final.pdf](https://tuhat.helsinki.fi/ws/portalfiles/portal/167373071/_harrison_mfula_ccaaS_final.pdf)
- [4]. L. Osmani, T. Kauppinen, M. Komu, and S. Tarkoma, "Multi-Cloud Connectivity for Kubernetes in 5G Networks," IEEE Communications Magazine, vol. 59, no. 10, pp. 42–47, Oct. 2021. doi: <https://doi.org/10.1109/mcom.110.2100124>.
- [5]. J. Shah and D. Dubaria, "Building Modern Clouds: Using Docker, Kubernetes & Google Cloud Platform," IEEE Xplore, Jan. 01, 2019. doi: <https://doi.org/10.1109/CCWC.2019.8666479>. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/8666479>. [Accessed: Oct. 26, 2020].
- [6]. O. Tomarchio, D. Calcaterra, and G. D. Modica, "Cloud Resource Orchestration in the Multi-Cloud Landscape: A Systematic Review of Existing Frameworks," Journal of Cloud Computing, vol. 9, no. 1, Sep. 2020. doi: <https://doi.org/10.1186/s13677-020-00194-7>.
- [7]. P. Suppala, "FinOps in SaaS Platform within Hybrid, Multi-Cloud, Multi-Tenant, Multi-Region Environments," Urn.fi, 2022. doi: <http://www.theseus.fi/handle/10024/780265>. Available: <https://urn.fi/URN:NBN:fi:amk-2022100420823>.

