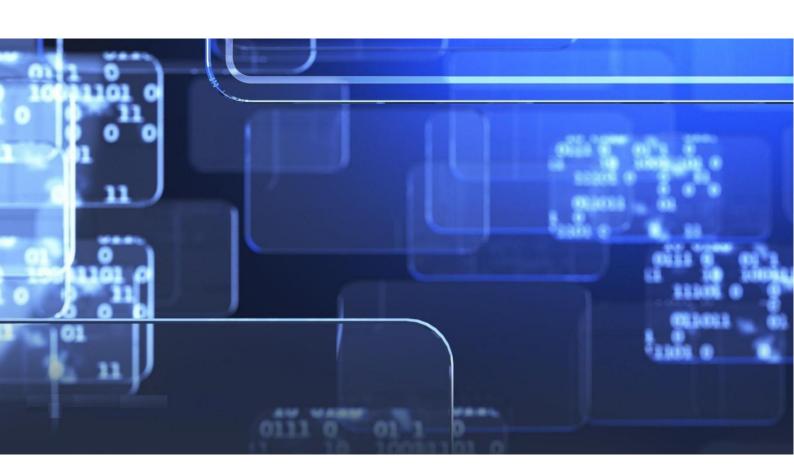


# KUBERNETES DEPLOYMENTS AND SCALING FUNDAMENTALS

**Training Material** 



# **Table of Contents**

Chapter 1: Introduction to Kubernetes Deployments	3
1.1 Overview of Kubernetes and its components	3
1.2 Understanding Deployments in Kubernetes	3
1.3 Difference between Deployments, ReplicaSets, and Pods	4
1.4 Benefits of Deployments in production environments	4
Chapter 2: Writing and Configuring Deployments	5
2.1 YAML basics for Kubernetes configurations	5
2.2 Defining a Deployment resource	5
2.3 Specifying replicas and container images	6
2.4 Adding environment variables using ConfigMaps and Secrets	6
2.5 Setting up resource limits and requests (CPU & memory)	8
Chapter 3: Scaling Applications in Kubernetes	8
3.1 Understanding horizontal and vertical scaling	8
3.2 Configuring horizontal scaling (HPA) with kubectl and YAML	9
3.3 Autoscaling based on CPU and custom metrics	9
3.4 Best practices and limitations of scaling	10

### **Chapter 1: Introduction to Kubernetes Deployments**

### 1.1 Overview of Kubernetes and Its Components

Kubernetes, often referred to as K8s, is a powerful open-source platform that automates the deployment, scaling, and management of containerized applications. Originally developed by Google, Kubernetes has become the standard for orchestrating containerized applications across clusters of machines.

### **Key Components of Kubernetes:**

### 1. Master Node:

- Role: Acts as the control plane and manages the cluster.
- c Components:
  - API Server: Exposes the Kubernetes API and acts as a gateway for communication between components.
  - Controller Manager: Governs controllers that regulate the state of the cluster, ensuring that the desired state matches the actual state.
  - **Scheduler**: Assigns workloads (Pods) to worker nodes based on resource availability and scheduling policies.
  - **Etcd**: A distributed key-value store that holds the configuration data and state of the cluster.

### 2. Worker Nodes:

- **Role**: Host the containers running the applications.
- components:
  - **Kubelet**: An agent that communicates with the master node and ensures that the containers in the Pods are running as expected.
  - Container Runtime: Software responsible for running containers,
     such as Docker or containerd.
  - **Kube Proxy**: Manages network communication for the Pods, facilitating load balancing and service discovery.
- 3. **Pods**: The smallest deployable units in Kubernetes, which can contain one or more containers. They share storage and network resources, making communication between containers efficient.

- 4. **Services**: Abstractions that define a logical set of Pods and a policy by which to access them, allowing for stable networking.
- 5. **Volumes**: Persistent storage resources that can be used by Pods, ensuring data persistence beyond the lifecycle of individual Pods.

Understanding these components is crucial for grasping how Kubernetes orchestrates and manages containerized applications, particularly when deploying and scaling applications.

# 1.2 Understanding Deployments in Kubernetes

A Deployment is a Kubernetes resource that provides declarative updates to applications. It allows developers to define the desired state for their applications, and Kubernetes takes care of maintaining that state over time.

### **Key Features of Deployments:**

- **Declarative Configuration**: Users specify the desired state in a YAML or JSON file, allowing Kubernetes to automatically adjust the actual state to match.
- **Rollouts**: Support for gradual rollout of changes to ensure minimal downtime and high availability during updates.
- **Scaling**: Easy to scale applications up or down by changing the number of replicas.
- **Health Checks**: Integrates liveness and readiness probes to monitor application health and availability.

### **Use Cases:**

- Continuous Deployment: Automate the deployment of applications after successful builds.
- Version Control: Manage different versions of applications with ease, enabling quick rollbacks.

# 1.3 Difference Between Deployments, ReplicaSets, and Pods

# **Components Explained:**

### 1. **Pods**:

The basic execution unit in Kubernetes, representing a single instance of a running process. Pods can encapsulate one or more containers.

### 2. **ReplicaSets**:

Ensures that a specified number of pod replicas are running at any time. If
a Pod fails or is deleted, the ReplicaSet automatically creates a new one to
maintain the desired count. ReplicaSets are usually managed by
Deployments.

### 3. **Deployments**:

 Higher-level constructs that manage ReplicaSets and provide declarative updates to Pods. They simplify the process of managing application updates, scaling, and rollback.

# **Relationship Summary:**

• **Pods** are the smallest units, **ReplicaSets** ensure availability, and **Deployments** provide an abstraction for managing and updating those ReplicaSets seamlessly.

### 1.4 Benefits of Deployments in Production Environments

Deployments offer several advantages, particularly in production scenarios:

- **Rolling Updates**: Deployments support rolling updates, allowing for new versions of applications to be deployed without taking the entire service offline.
- **Automatic Rollbacks**: If a deployment fails, Kubernetes can revert to the previous stable version automatically.
- **Declarative Management**: Users define their desired state, enabling Kubernetes to manage the lifecycle of the application.
- **Horizontal Scaling**: Deployments facilitate scaling up or down by simply changing the number of replicas.

- **High Availability**: By maintaining multiple replicas and performing health checks, Deployments enhance the availability of applications.
- Simplified Management: Deployments abstract the complexity of managing Pods and ReplicaSets, providing a more straightforward interface for application management.

### **Chapter 2: Writing and Configuring Deployments**

# 2.1 YAML Basics for Kubernetes Configurations

YAML (YAML Ain't Markup Language) is a human-readable data serialization standard commonly used in configuration files. It is the preferred format for Kubernetes manifests.

# **Key Characteristics of YAML:**

- **Indentation**: Structure is defined through indentation; consistency is key.
- **Key-Value Pairs**: Data is organized as key-value pairs, with the key being a string and the value being a string, number, or list.
- **Lists**: Denoted with a hyphen (-), allowing for the representation of multiple items.

### **Example of YAML Structure:**

```
apiVersion: v1
kind: Service
metadata:
name: my-service
spec:
ports:
- port: 80
targetPort: 8080
selector:
app: my-app
```

Understanding YAML syntax is essential for defining Kubernetes resources accurately.

### 2.2 Defining a Deployment Resource

To create a Deployment in Kubernetes, you must define it using a YAML file. The structure includes specifying the API version, kind, metadata, and specifications for the Pods.

### **Example Deployment Definition:**

```
apiVersion: apps/v1
kind: Deployment
metadata:
 name: my-app-deployment
spec:
 replicas: 3
 selector:
  matchLabels:
   app: my-app
 template:
  metadata:
   labels:
    app: my-app
  spec:
   containers:
   - name: my-app-container
    image: my-app-image:latest
    ports:
    - containerPort: 8080
```

# **Explanation of Key Fields:**

- apiVersion: Specifies the API version used for the Deployment object.
- **kind**: Indicates the type of resource (Deployment).
- **metadata**: Contains metadata information such as the name and labels.
- spec: Details the desired state, including replicas and Pod templates.

2.3 Specifying Replicas and Container Images

In the Deployment specification, you can define the number of replicas and the container

image that will run in each Pod.

**Specifying Replicas:** 

The replicas field indicates how many identical Pods should be running. This helps ensure

high availability.

**Container Images:** 

The image field within the container specification indicates which Docker image to use. It

is a best practice to tag images with specific version tags rather than using latest to avoid

unintentional updates.

**Example:** 

spec:

replicas: 3

template:

spec:

containers:

- name: my-app-container

image: my-app-image:v1.0

2.4 Adding Environment Variables Using ConfigMaps and Secrets

Kubernetes allows the injection of environment variables into containers via ConfigMaps

for non-sensitive information and Secrets for sensitive information.

**Using ConfigMaps:** 

ConfigMaps are used to store non-sensitive configuration data. They can be referenced in

the Deployment to pass configuration values as environment variables.

8

# Example of ConfigMap:

apiVersion: v1

kind: ConfigMap

metadata:

name: my-config

data:

DATABASE\_URL: "mysql://user:password@hostname:port/dbname"

# **Using Secrets:**

Secrets are designed to store sensitive data, such as passwords or API keys, in a more secure way than ConfigMaps.

# Example of Secret:

apiVersion: v1

kind: Secret

metadata:

name: my-secret

type: Opaque

data:

PASSWORD: cGFzc3dvcmQ= # Base64 encoded

# **Referencing in Deployment:**

You can reference these in your Deployment manifest to make the values available to your application.

# Example Deployment with Environment Variables:

env:

- name: DATABASE\_URL

valueFrom:

configMapKeyRef:

name: my-config

key: DATABASE\_URL

- name: PASSWORD

valueFrom:

secretKeyRef:

name: my-secret key: PASSWORD

# 2.5 Setting Up Resource Limits and Requests (CPU & Memory)

Resource management is crucial in Kubernetes to ensure fair resource allocation and optimize application performance. You can define resource requests and limits in your Deployment configuration.

# **Requests vs. Limits:**

- **Requests**: The minimum amount of CPU/memory that Kubernetes guarantees for a container. Used for scheduling Pods.
- **Limits**: The maximum amount of CPU/memory that a container can use. If a container exceeds this limit, it may be throttled or killed.

### **Example of Resource Configuration:**

resources:

requests:

memory: "64Mi"

cpu: "250m"

limits:

memory: "128Mi"

cpu: "500m"

### **Best Practices:**

- Set requests and limits to appropriate levels based on application profiling.
- Monitor resource usage and adjust as necessary.

### **Chapter 3: Scaling Applications in Kubernetes**

## 3.1 Understanding Horizontal and Vertical Scaling

Scaling is a critical aspect of managing applications in Kubernetes, allowing them to handle varying loads efficiently.

# **Horizontal Scaling (Scaling Out/In):**

- **Definition**: Involves adding more Pods to handle increased load (scaling out) or removing Pods when demand decreases (scaling in).
- Use Case: Typically used for stateless applications, where multiple instances can handle requests simultaneously.

### **Vertical Scaling (Scaling Up/Down):**

- **Definition**: Involves increasing or decreasing the resources (CPU/memory) allocated to existing Pods.
- **Use Case**: Suitable for stateful applications where a single instance needs more resources to perform optimally.

### **Considerations:**

- Kubernetes is primarily designed for horizontal scaling due to its orchestration capabilities.
- Vertical scaling is limited by the resources available on the node.

# 3.2 Configuring Horizontal Scaling (HPA) with kubectl and YAML

Horizontal Pod Autoscaler (HPA) automatically adjusts the number of Pods in a Deployment based on observed CPU utilization or other select metrics.

### **Basic HPA Configuration:**

To create an HPA, you can use the kubectl command or define it in a YAML file.

# HPA Example via kubectl:

kubectl autoscale deployment my-app-deployment --cpu-percent=50 --min=2 --max=10

### **YAML Definition of HPA:**

apiVersion: autoscaling/v1

kind: HorizontalPodAutoscaler

metadata:

name: my-hpa

spec:

scaleTargetRef:

apiVersion: apps/v1

kind: Deployment

name: my-app-deployment

minReplicas: 2

maxReplicas: 10

targetCPUUtilizationPercentage: 50

### **Key Fields Explained:**

- scaleTargetRef: References the Deployment to be scaled.
- **minReplicas** and **maxReplicas**: Define the minimum and maximum number of replicas.
- **targetCPUUtilizationPercentage**: Specifies the target CPU utilization percentage that triggers scaling.

### 3.3 Autoscaling Based on CPU and Custom Metrics

Kubernetes supports scaling based on standard metrics like CPU, but also allows scaling based on custom metrics collected from the application.

### **Setting Up Custom Metrics:**

To enable custom metrics, you need a metrics server and can use tools like Prometheus to gather these metrics.

### Example of Custom Metric Configuration in HPA:

```
spec:
  metrics:
  - type: Object
  object:
    metric:
    name: requests_per_second
  target:
    type: AverageValue
    averageValue: 100
```

# **Explanation:**

This configuration will scale the Pods based on the average number of requests per second, providing a dynamic way to adjust to varying loads.

# 3.4 Best Practices and Limitations of Scaling

### **Best Practices:**

- **Monitor Resource Usage**: Use monitoring tools (like Prometheus and Grafana) to track application performance and resource utilization regularly.
- **Set Reasonable Limits**: Define clear resource limits and requests to prevent resource contention and ensure fair resource allocation among Pods.
- Use Readiness and Liveness Probes: Implement probes to ensure that Pods are healthy and ready to handle traffic, which helps during scaling operations.

### **Limitations:**

- Cluster Resources: The ability to scale is limited by the available resources in the Kubernetes cluster. Ensure you have sufficient resources allocated to support scaling.
- **Scaling Lag**: Autoscaling may introduce delays in scaling actions. Proper thresholds and metrics should be set to minimize these delays.

• **Complexity of Custom Metrics**: Implementing autoscaling based on custom metrics may require additional tooling and setup, increasing complexity.

### Conclusion

This comprehensive guide covers the essential aspects of Kubernetes Deployments, from understanding the core components to writing configurations and managing scalability. By following best practices and leveraging Kubernetes' powerful features, you can ensure that your applications are resilient, scalable, and maintainable in a cloud-native environment.