# Autoscaling in Kubernetes

Abdullah Al Hasib

abdullah.alhasib@kratosdefense.com

09-02-2022

**1** Kubernetes Recap

**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

**1** Kubernetes Recap

**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

## A Quick Recap

**What is kubernetes:**

# A Quick Recap

**Kubernetes architecture:**

**1** Kubernetes Recap

**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

## Autoscaling

**Autoscaling:**

- automatically adjusts computational resource usage according to the load.

**Why do we need it:**

- to cope with the demand
- to reduce cost
- to reduce power consumption

## Autoscaling Types

The Kubernetes autoscaling mechanism uses two layers:

- Pod-based scaling - supported by
    - Horizontal Pod Autoscaler (HPA)
    - Vertical Pod Autoscaler (VPA)
- Node-based scaling - supported by
    - Cluster Autoscaler (CA)

|            | Pods                | Nodes               |
|------------|---------------------|---------------------|
| Horizontal | # of pods           | # of nodes          |
| Vertical   | resources of a pod  | resources of a node |

**1** Kubernetes Recap
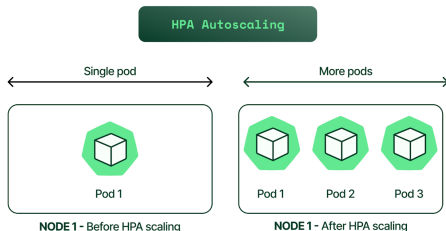
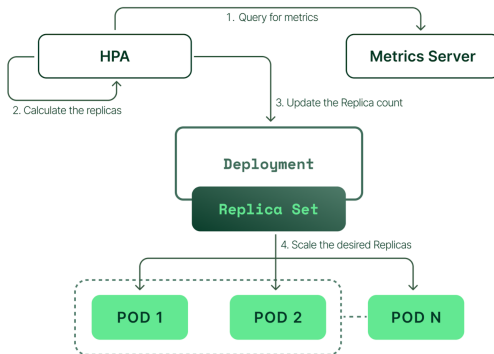**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

## Horizontal Pod Autoscaling

HPA increases or decreases the number of pods in a replication controller, deployment, replica set etc. based on CPU utilization

## How Does HPA Work?



```
desiredReplicas = ceil[currentReplicas * ( currentMetricValue / desiredMetricValue )]
```

## Metrics Types

Types of metrics APIs:

1. Resource Metrics (`metrics.k8s.io`)

 - predefined resource usage metrics of Pods and Nodes
 - can be expressed as raw value of percentage
 - example: CPU, Memory

Metrics Types

Types of metrics APIs:

1. Resource Metrics (`metrics.k8s.io`)
   - predefined resource usage metrics of Pods and Nodes
   - can be expressed as raw value of percentage
   - example: `CPU`, `Memory`
2. Custom Metrics (`custom.metrics.k8s.io`)
   - custom metrics associated with a Kubernetes object
   - example: rate_of_client_requests

## Metrics Types

Types of metrics APIs:

1. Resource Metrics (`metrics.k8s.io`)
   - predefined resource usage metrics of Pods and Nodes
   - can be expressed as raw value of percentage
   - example: CPU, Memory
2. Custom Metrics (`custom.metrics.k8s.io`)
   - custom metrics associated with a Kubernetes object
   - example: rate_of_client_requests
3. External Metrics (`external.metrics.k8s.io`)
   - custom metrics not associated with a Kubernetes object

Customizing Scaling Behavior

1. scaleup: control scaling behavior while scaling up
   - stabilizationWindowSeconds: (default:0)
   - selectPolicy: can be Min, Max, Disabled (default:Max)
   - policies
     - type: Pods or Percent
     - periodSeconds: (default: 60)
     - value

Customizing Scaling Behavior

**1** scaleup: control scaling behavior while scaling up
  - stabilizationWindowSeconds: (default:0)
  - selectPolicy: can be Min, Max, Disabled (default:Max)
  - policies
    - type: Pods or Percent
    - periodSeconds: (default: 60)
    - value

**2** scaledown: control scaling behavior while scaling down
  - stabilizationWindowSeconds: (default: 300)
  - selectPolicy: can be Min, Max, Disabled (default:Max)
  - policies
    - type: Pods or Percent
    - periodSeconds: (default: 60)
    - value

## User Stories

1. Scale up as fast as possible, scale down as usual

```
1  behavior:
2    scaleUp:
3      policies:
4      - type: Percent
5        value: 900%
```

## User Stories

**1** Scale up as fast as possible, scale down as usual

```
1  behavior:
2    scaleUp:
3      policies:
4      - type: Percent
5        value: 900%
```

**2** Scale up as usual, do not scale down

```
1  behavior:
2    scaleDown:
3      selectPolicy: Disabled
```

## User Stories

4. Scale Up As Fast As Possible, Scale Down Very Gradually

```
 1  behavior:
 2    scaleUp:
 3      policies:
 4      - type: Percent
 5        value: 900%
 6    scaleDown:
 7      policies:
 8      - type: Pods
 9        value: 1
10        periodSeconds: 600
```

## User Stories (Cont'd)

**5** Stabilization before scaling down

```
1  behavior:
2    scaleDown:
3      stabilizationWindowSeconds: 600
4      policies:
5      - type: Pods
6        value: 5
```

# HPA Demo

**1** Kubernetes Recap
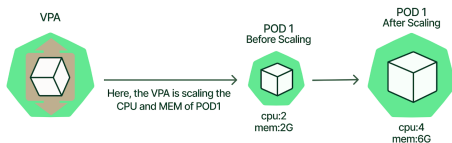
**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

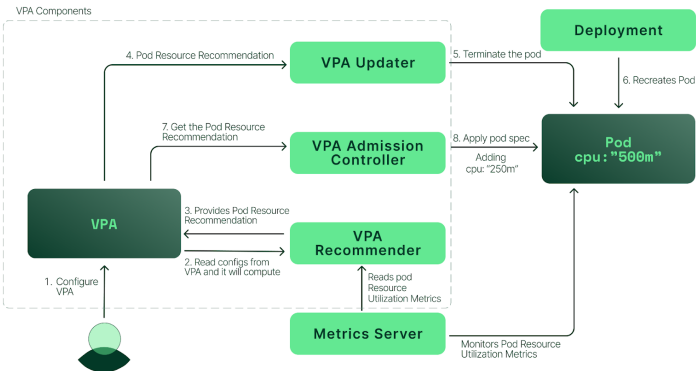**6** Conclusion

## Vertical Pod Autoscaling

VPA adjusts the resource requests and limits of containers in the cluster.



Resource configuration types:

- requests: define the minimum amount of resources that containers need
- limits: define the maximum amount of resources that a given container can consume

# How Does VPA Work?

## VPA Concepts

Types Operation Modes:

- Off: provides recommendations only
- Initial: only assigns resource requests on pod creation
- Recreate: assigns resource requests on pod creation and updates them on existing pods by evicting them. [should be used rarely]
- Auto: currently does the same as Recreate. [may cause service downtime]

# VPA Demo

**1** Kubernetes Recap
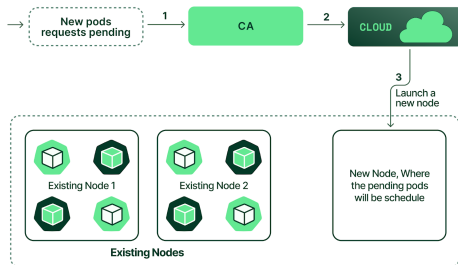
**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

## Cluster Autoscaling

CA adjusts the number of nodes in the cluster when pods fail to schedule or when nodes are underutilized

**1** Kubernetes Recap

**2** Autoscaling

**3** HPA

**4** VPA

**5** CA

**6** Conclusion

## Conclusion

- Install a metric server
- Do not mix HPA with VPA
- Define pod requests and limits
- Resource requests should be close to the average usage of the pods
- Increase CPU limits for slow starting applications

[1] Minikube Handbook
[2] Horizontal Pod Autoscaling
[3] Vertical Pod Autoscaling
[4] Scaling Kubernetes Clusters
[5] KEDA: Kubernetes Event-Driven Autoscaling

*Thanks!*