# Autoscaling in Kubernetes

Abdullah Al Hasib

abdullah.alhasib@kratosdefense.com

09-02-2022

**1** Introduction

**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

**5** Cluster Autoscaling

**6** Conclusion

**1** Introduction

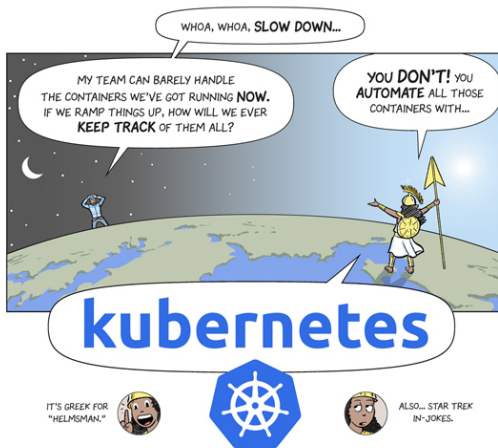**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

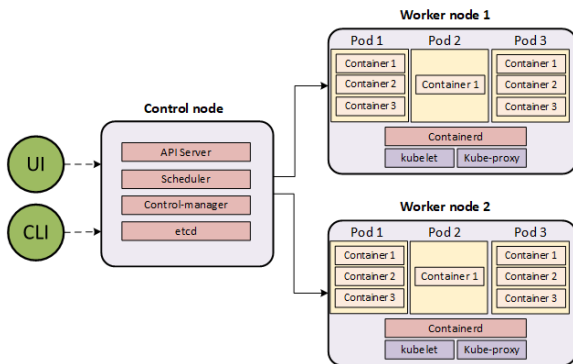**5** Cluster Autoscaling

**6** Conclusion

## A Quick Recap

**What is kubernetes:**

## A Quick Recap

**Kubernetes architecture:**

**1** Introduction

**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

**5** Cluster Autoscaling

**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)
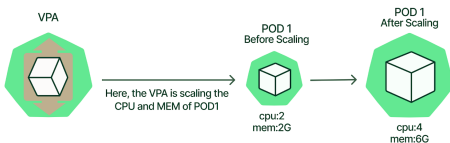- Node-based scaling - supported by
    - Cluster Autoscaler (CA)

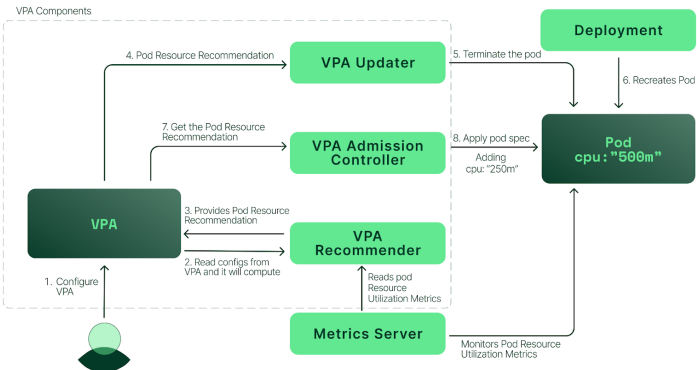|            | Pods               | Nodes               |
|------------|--------------------|---------------------|
| Horizontal | # of pods          | # of nodes          |
| Vertical   | resources of a pod | resources of a node |

## Virtical 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

Introduction
○○○

Motivation
○○○

**Vertical Pod Autoscaling**
○○●○○○○

Horizontal Pod Autoscaling
○○○○○

Cluster Autoscaling
○○○

Conclusion
○○○○

# How Does VPA Work?

## VPA Concepts

Types Operation Modes:

- Off: default mode
- Auto: may cause service downtime
- Recreate: should be used rarely
- Initial

Types of metrics server APIs:

- Resource Metrics: Predefined resource usage metrics (CPU and memory)
- Custom Metrics: Custom metrics associated with a K8s object
- External Metrics: Custom metrics not associated with a K8s object

## Limitations

- Do not use VPA with HPA while they are using the same resource metrics such as CPU and memory usage.
- VPA might recommend more resources than available in the cluster, thus causing the pod to not be assigned to a node.
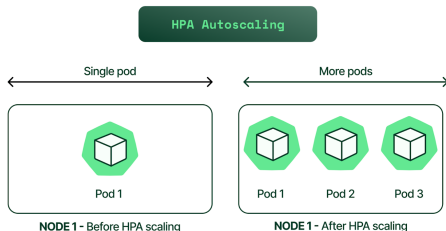- To overcome this limitation, set the LimitRange to the maximum available resources.

VPA Demo

# VPA Demo

**1** Introduction

**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

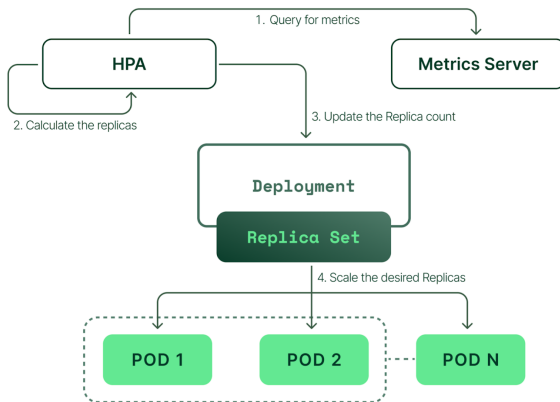**5** Cluster Autoscaling

**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?

## Limitations

- One of HPA's most well-known limitations is that it does not work with DaemonSets.
- If you don't efficiently set CPU and memory limits on pods, your pods may terminate frequently or, on the other end of the spectrum, you'll waste resources.
- If the cluster is out of capacity, HPA can't scale up until new nodes are added to the cluster. Cluster Autoscaler (CA) can automate this process.
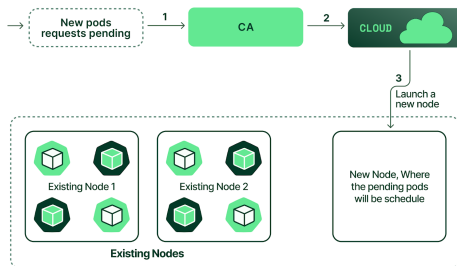
Introduction  Motivation  Vertical Pod Autoscaling  **Horizontal Pod Autoscaling**  Cluster Autoscaling  Conclusion
000        000        000000                   00000                              000                 0000

HPA Demo

# HPA Demo

**1** Introduction

**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

**5** Cluster Autoscaling

**6** Conclusion

## Cluster Autoscaling

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

## Limitations

- CA does not make scaling decisions using CPU or memory usage. It only checks a pod's requests and limits for CPU and memory resources. This limitation means that the unused computing resources requested by users will not be detected by CA, resulting in a cluster with waste and low utilization efficiency.

- Whenever there is a request to scale up the cluster, CA issues a scale-up request to a cloud provider within 30–60 seconds. The actual time the cloud provider takes to create a node can be several minutes or more. This delay means that your application performance may be degraded while waiting for the extended cluster capacity.

**1** Introduction

**2** Motivation

**3** Vertical Pod Autoscaling

**4** Horizontal Pod Autoscaling

**5** Cluster Autoscaling

**6** Conclusion

## Conclusion

- Here are two best practices for making efficient use of Cluster Autoscaler: Ensure resource availability for the Cluster Autoscaler pod—you can do that by defining a minimum of one CPU for resource requests made to the cluster autoscaler pod. It is critical to ensure that the node running the cluster autoscaler pod has enough resources. Otherwise, the cluster autoscaler might become non responsive.

- Ensure all pods have a defined resource request—to function correctly, the cluster autoscaler needs a specified resource request. This is because the cluster autoscaler makes decisions based on the status of pods and utilization of individual nodes, and can be thrown off if the calculation is off.

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

*Thanks!*