



kubernetes



Intro to Kubernetes Autoscaling

HPA, VPA, CA and beyond
William Albertus Dembo



CLOUD NATIVE
COMPUTING FOUNDATION

Outline

Introduction

Problems
K8s Components

01

02

HPA

How it work?
Demo
Problem

03

CA

How it work?
Problem

04

VPA

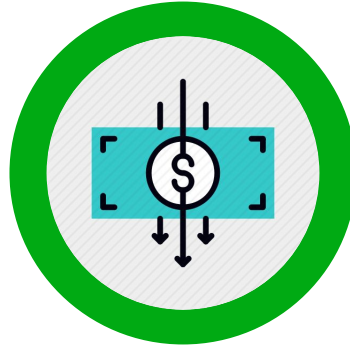
How it work?
Demo

Problems



Traffic Spike

Lorem Ipsum is simply dummy text of the printing and typesetting industry.



Cost saving

Lorem Ipsum is simply dummy text of the printing and typesetting industry.

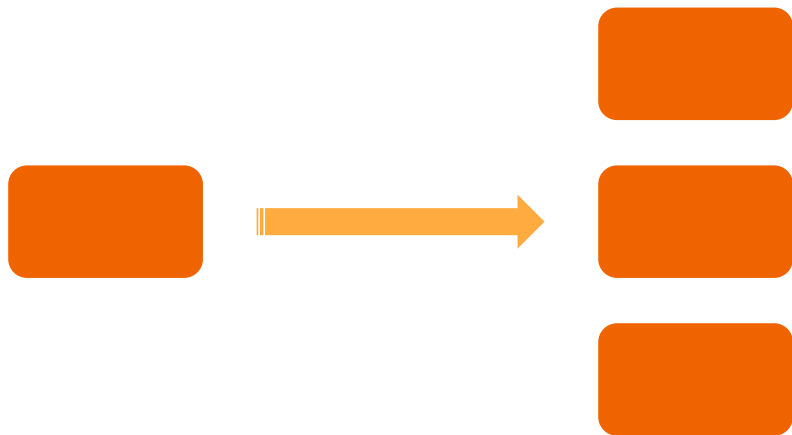


Maintenance

Lorem Ipsum is simply dummy text of the printing and typesetting industry.

Scaling Dimensions

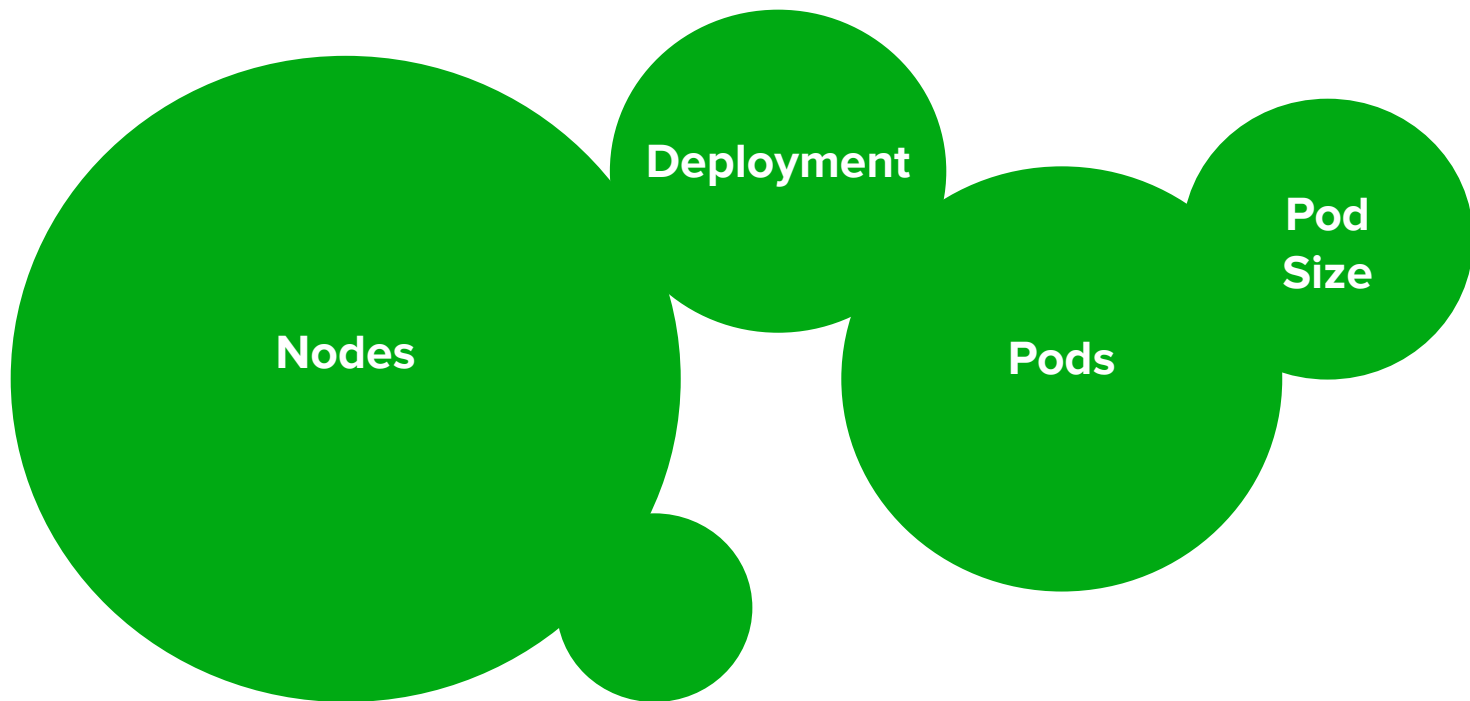
Horizontal Scaling



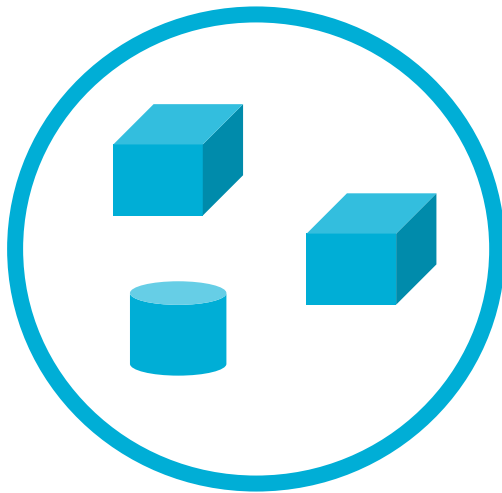
Vertical Scaling



Kubernetes Components

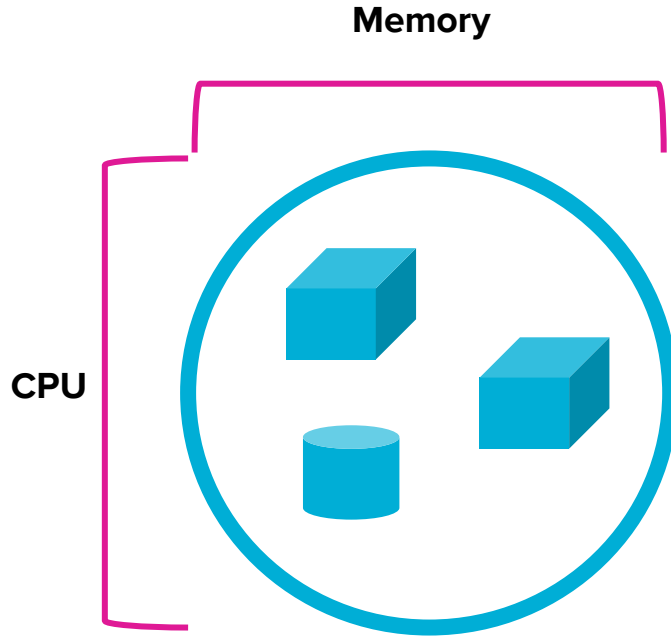


Pod



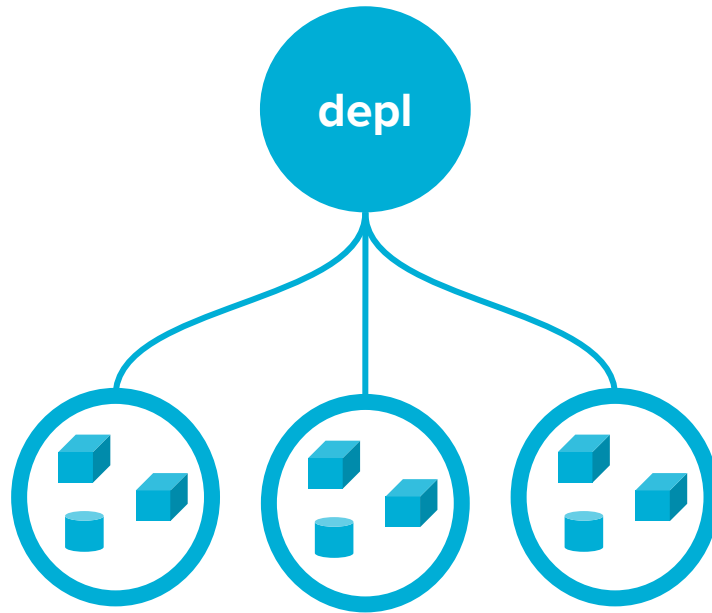
**A Pod: Group of containers
deployed together on
same host**

Pod Size



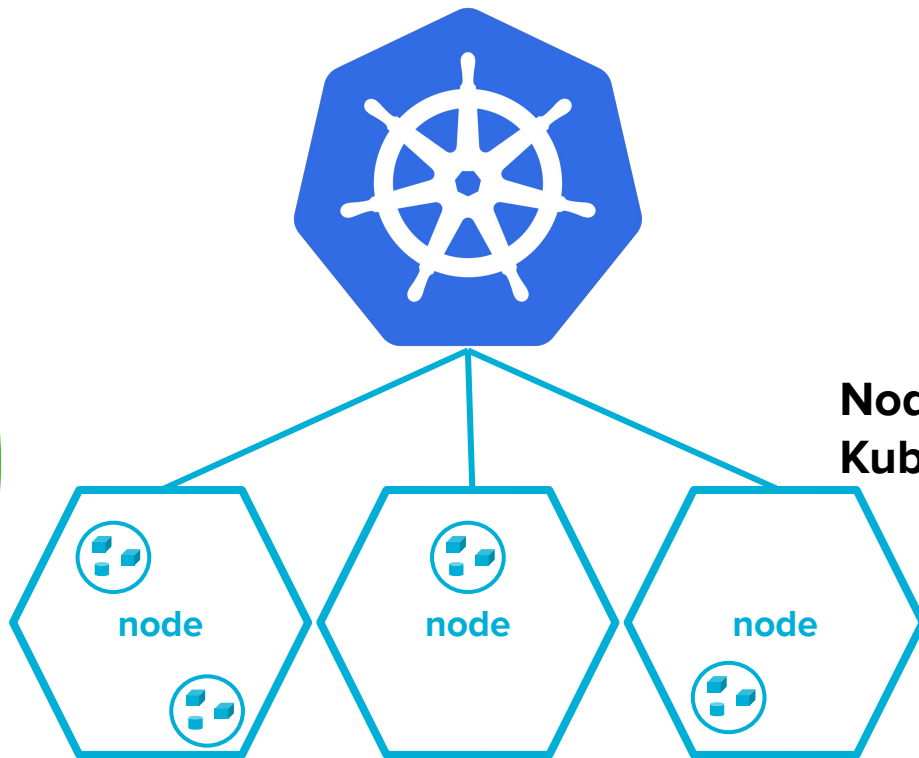
Pod size: Request and Limit assigned to Pod (Not actual usage)

Deployment



**Deployment: Maintains
homogenous set of pods**

node



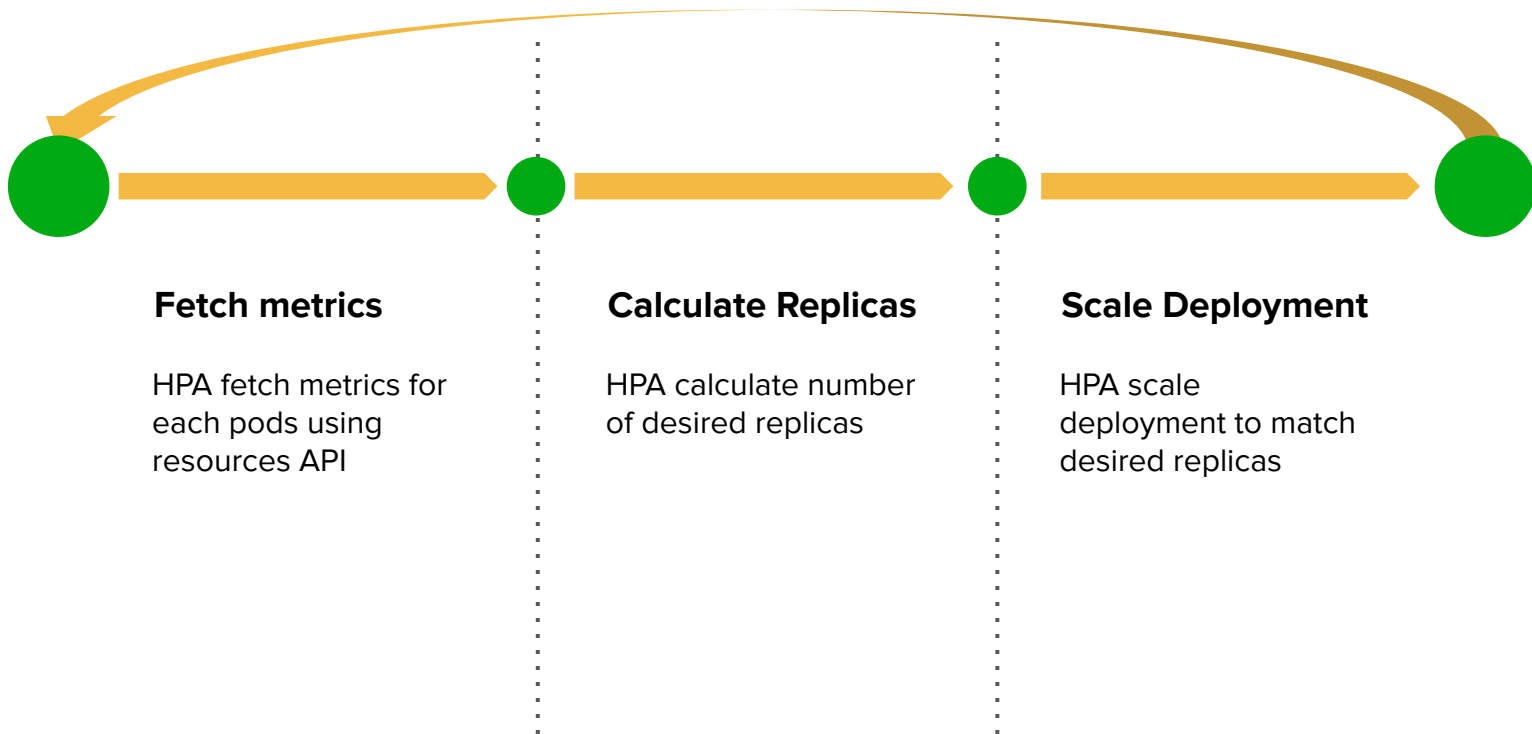
**Node: Worker machine in
Kubernetes cluster**

Kubernetes Scaling Dimensions

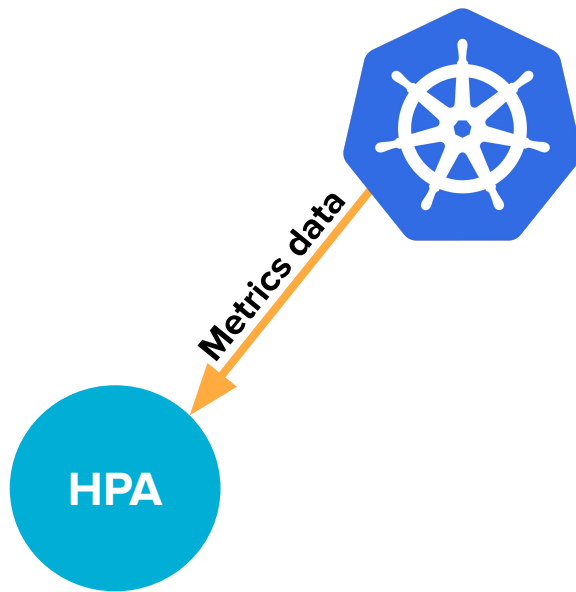
	Nodes	Pods
Horizontal	Number of nodes	Number of pods
Vertical	Size of node	Size of pod

Horizontal Pod Autoscaler

How it work?



HPA



HPA

Desired
Replicas

HPA

Metrics data



HPA

Desired
Replicas

HPA

Scale

depl

Metrics data



HPA

Desired
Replicas

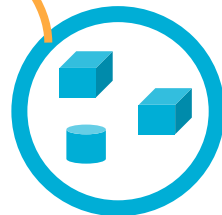
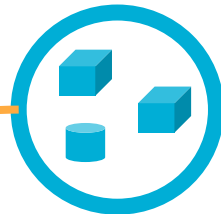
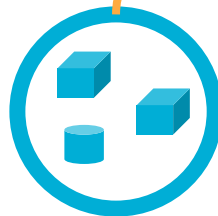
HPA

Metrics data

Scale

depl

Spawn/Kill



HPA

Desired
Replicas

HPA

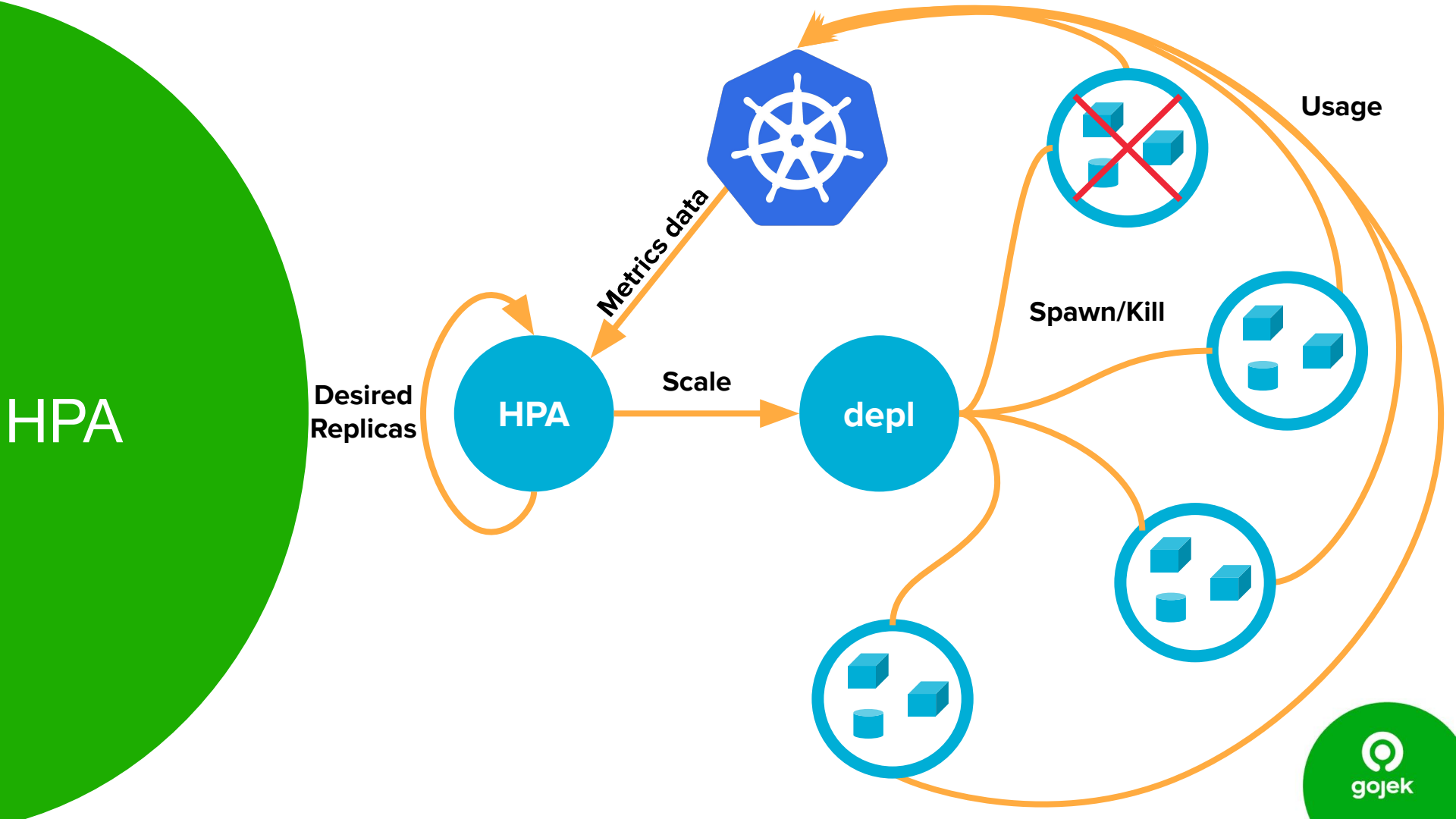
Scale

depl

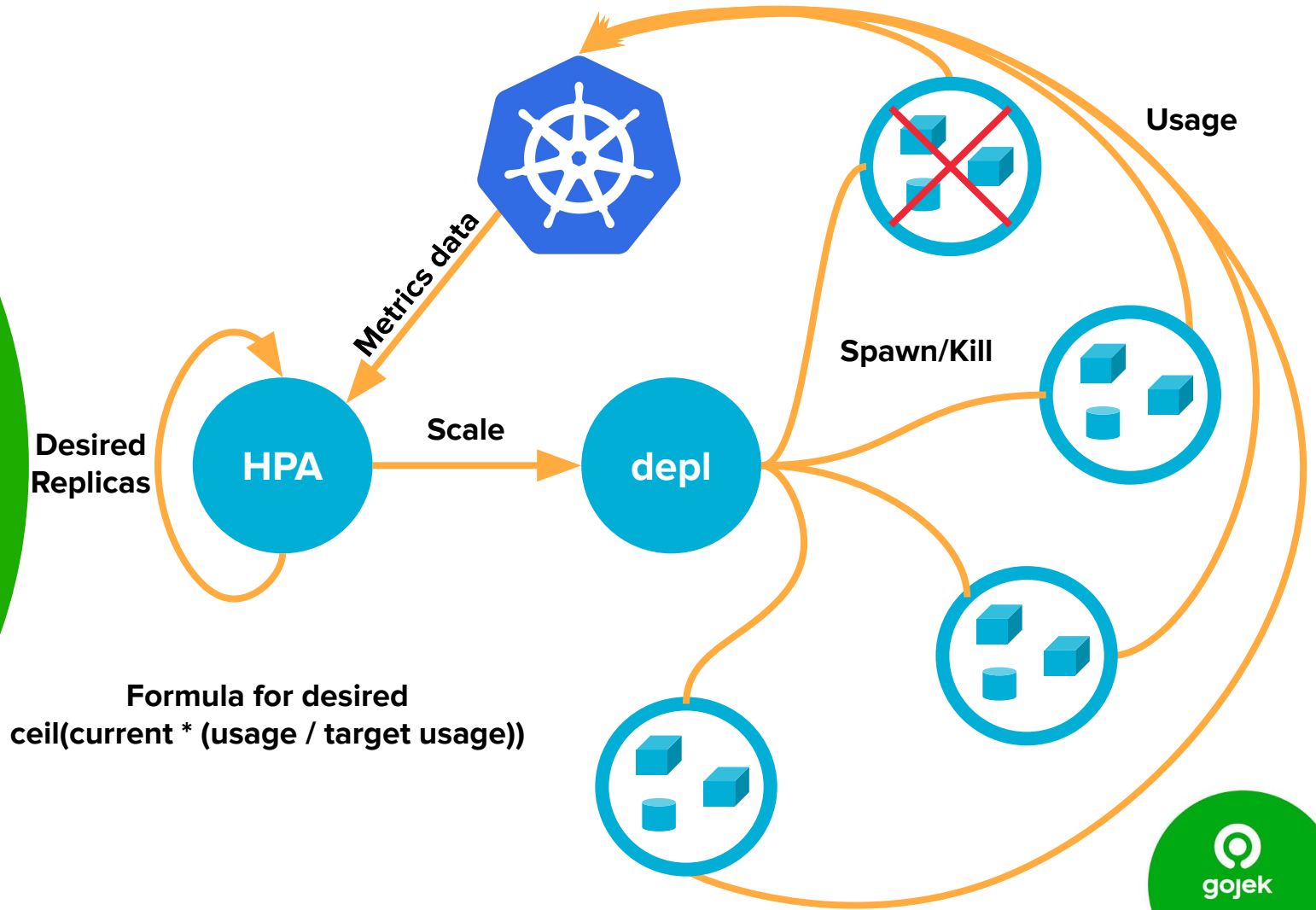
Metrics data

Spawn/Kill

Usage

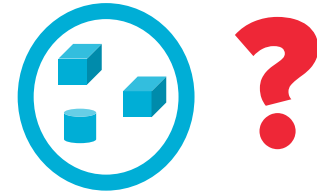
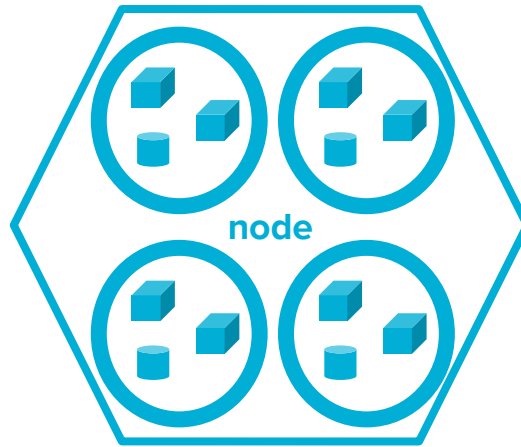
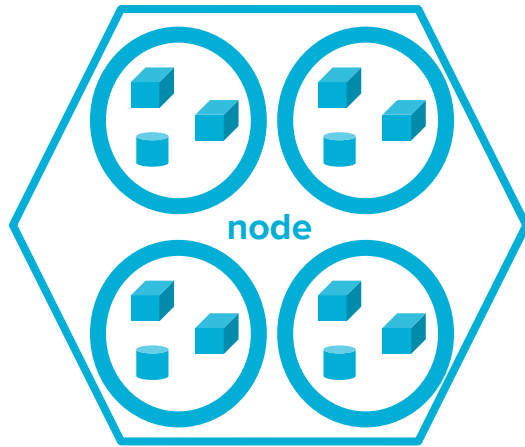


HPA



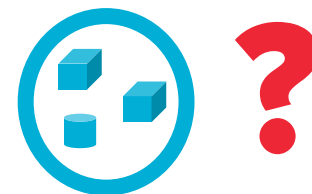
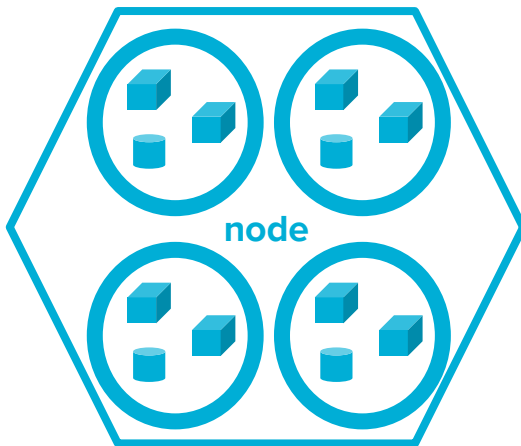
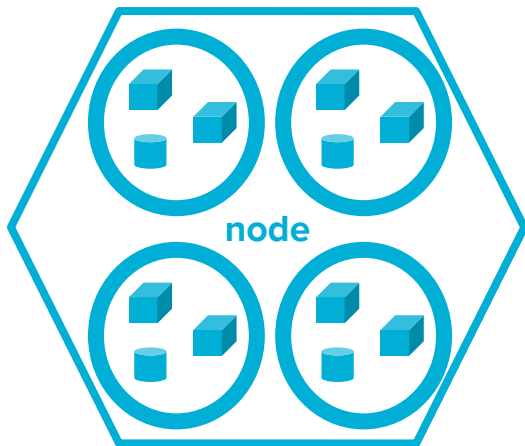


HPA Problem

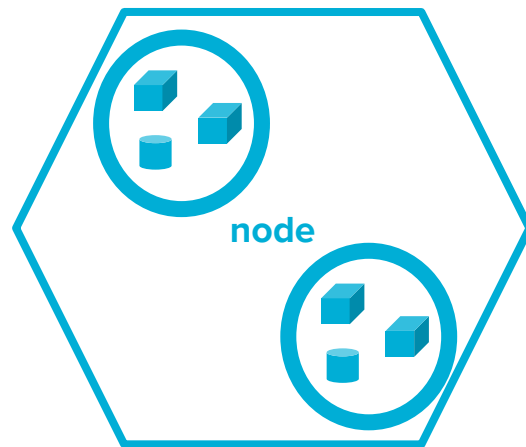
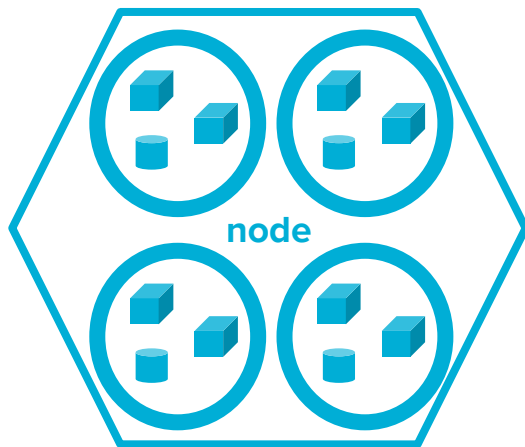
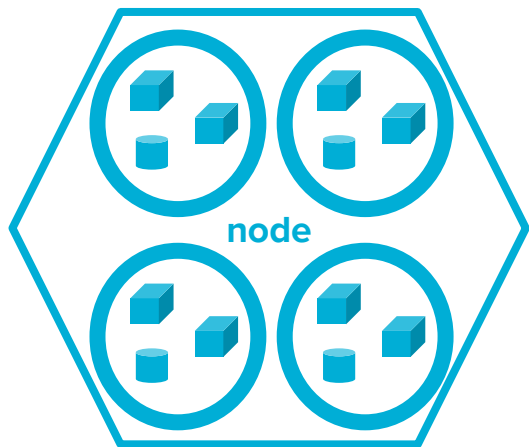


Cluster Autoscaler

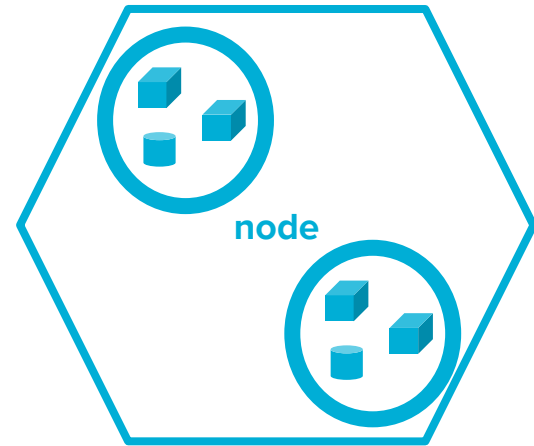
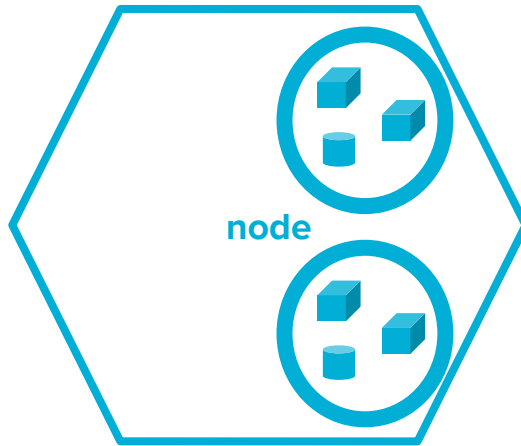
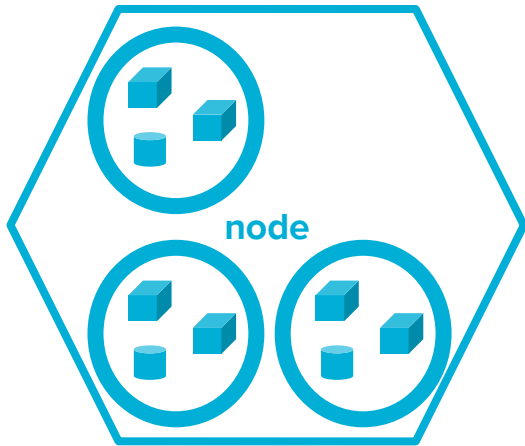
Solution?



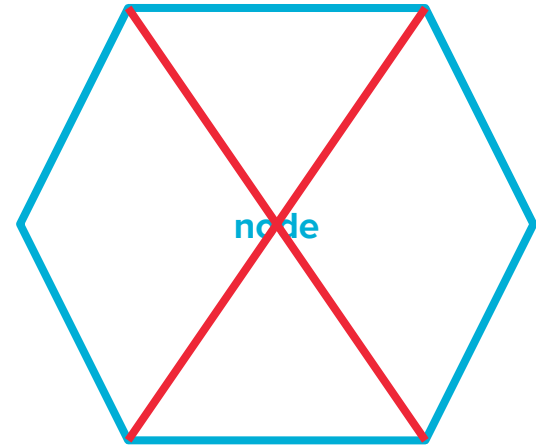
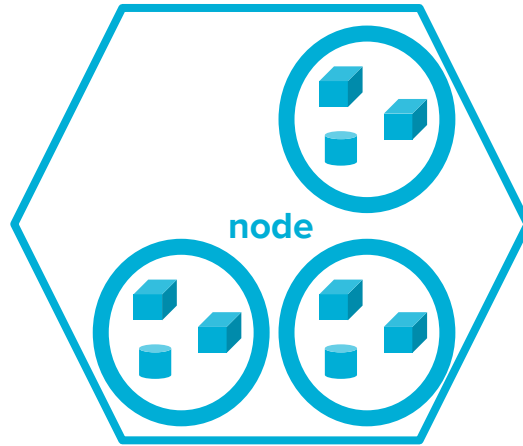
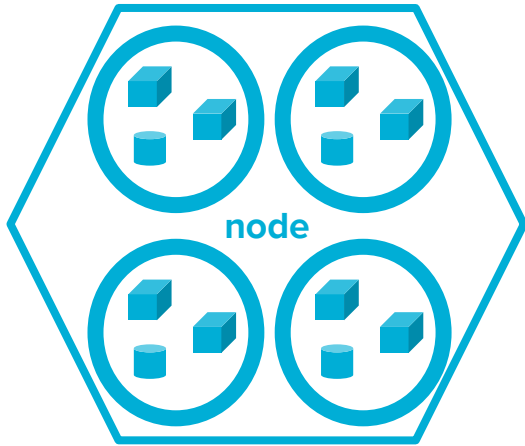
Spawn more node!



Underutilized resource?



Delete node



CA rules

Spawn node when there are pods that are unable to schedule due to insufficient resource

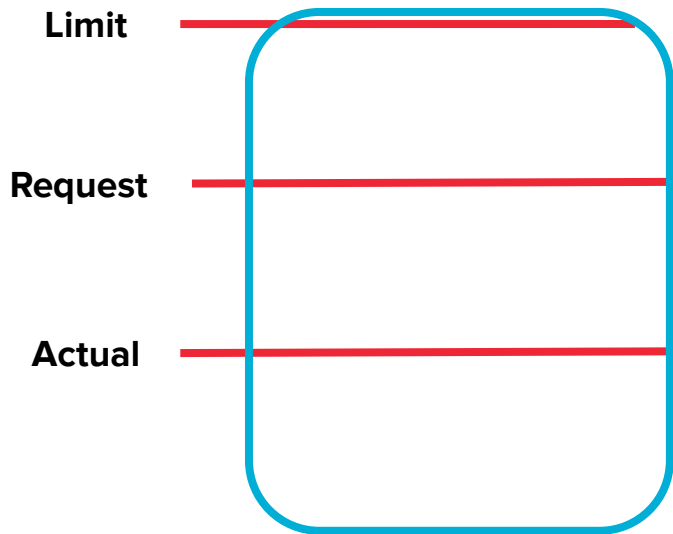
**Kill node when some nodes are consistently unneeded for a significant amount of time.
Unneeded nodes mean that it has low utilization and important pods are able to relocate to other nodes**

DEMO?

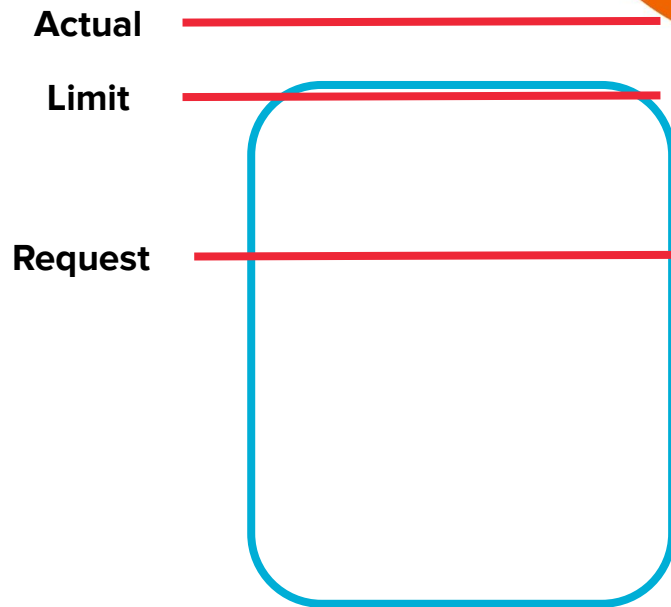
CA resources

- <https://github.com/kubernetes/autoscaler/tree/master/cluster-autoscaler>
- <https://aws.amazon.com/premiumsupport/knowledge-center/eks-cluster-autoscaler-setup/>
-

Utilization Problem



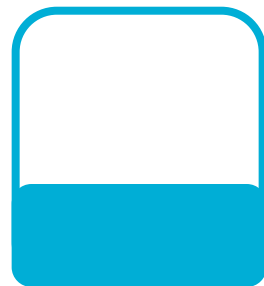
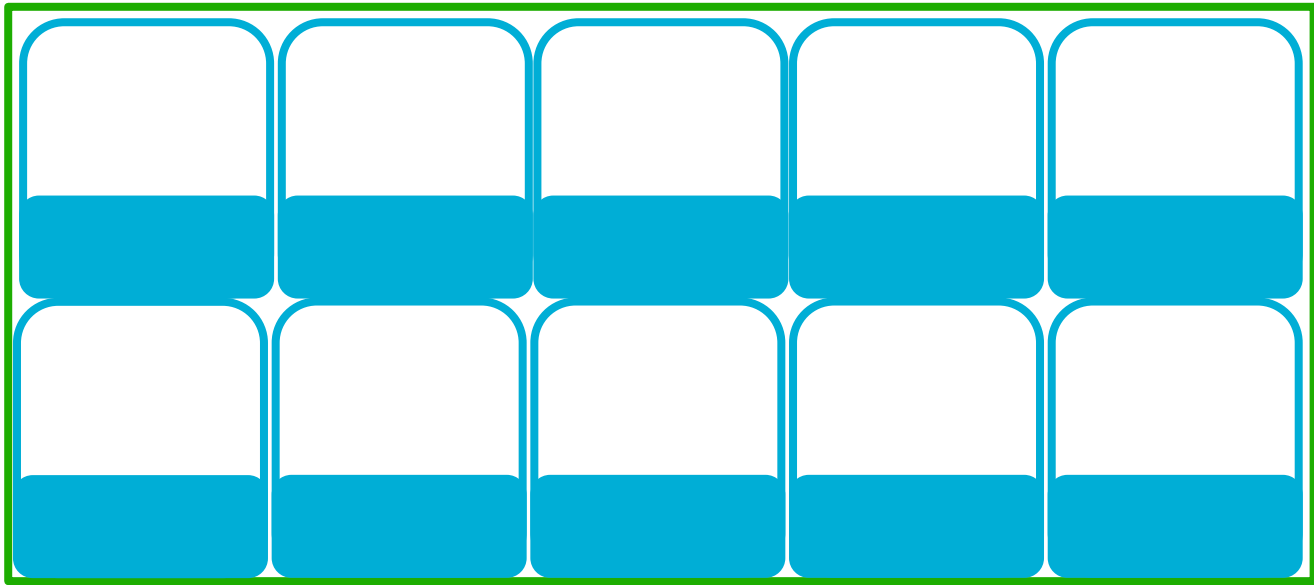
Too much



Not enough

Size does matter

Node



Vertical Pod Autoscaler

VPA Components



Recommender

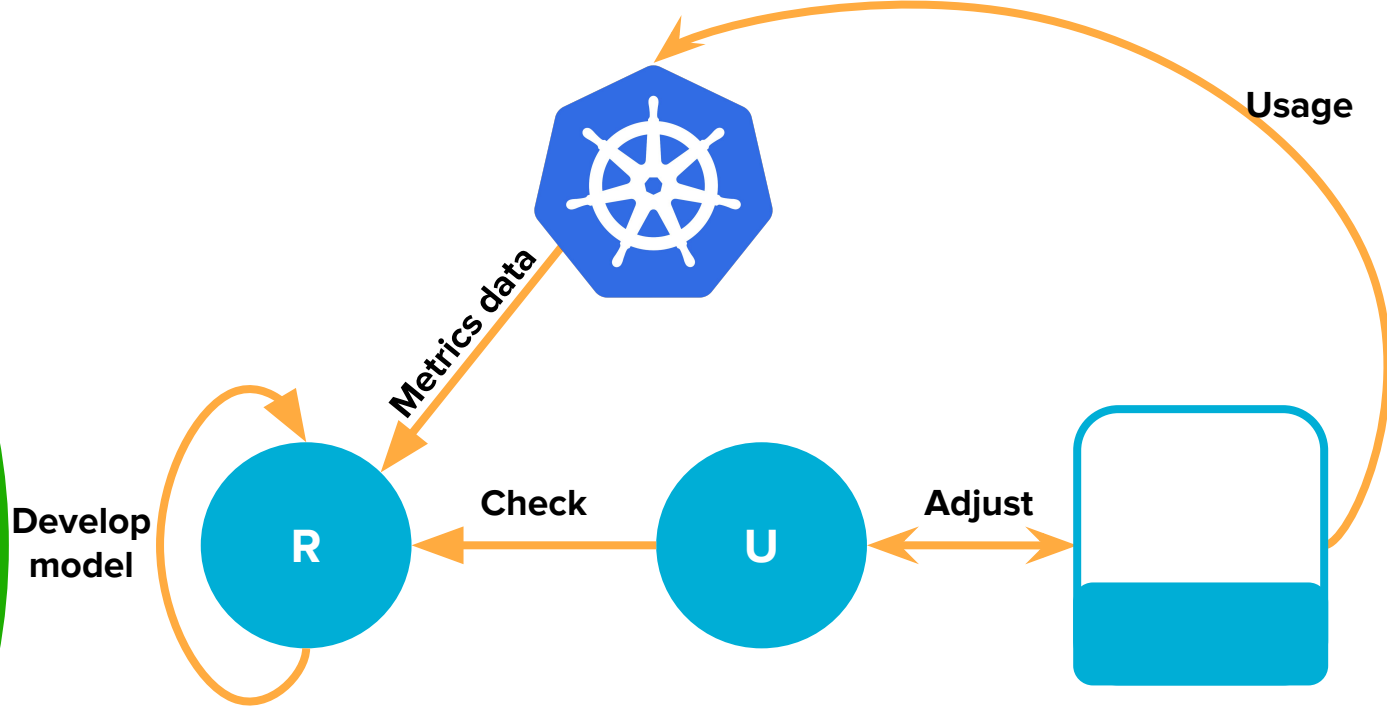
Monitors resources consumption and provides recommended values.



Updater

Check which pods doesn't have correct resources set and update them.

VPA



Result

Node

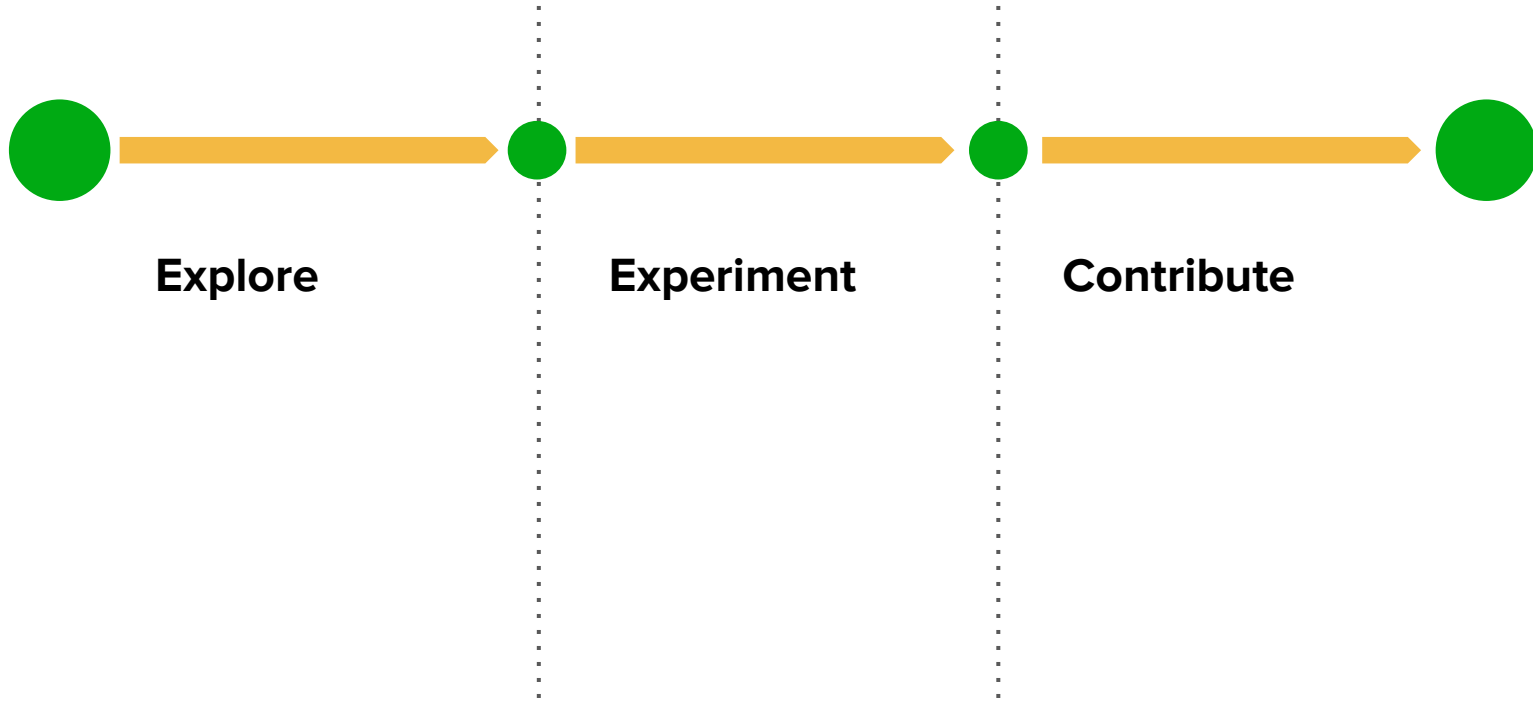




Limitations

- Cannot use HPA and VPA together
- VPA is beta (some features might introduce downtime)
- HPA stable only support CPU (v2beta2 support memory and custom)
- CA depends on Cloud Provider
-

Beyond



Thank You!



FAQ

#

Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s

#

Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s

#

Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s

#

Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s