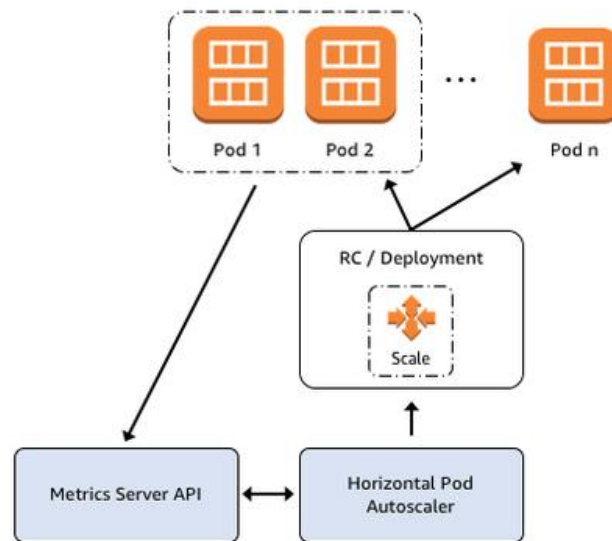# Horizontal Pod Autoscaling

The Horizontal Pod Autoscaler changes the shape of your Kubernetes workload by automatically increasing or decreasing the number of Pods in response to the workload's CPU or memory consumption, or in response to custom metrics reported from within Kubernetes or external metrics from sources outside of your cluster
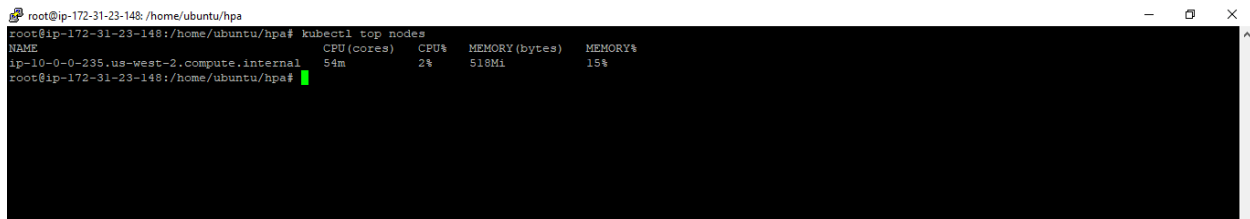
Architectural Diagram:



As shown in the image below, until and unless metric server is not deployed in the cluster, it will not show the real time metrics of the system.



Install Metrics Server commands:

– kubectl apply -f https://github.com/kubernetes-sigs/metrics-server/releases/latest/download/components.yaml

– kubectl get deployment metrics-server -n kube-system

As shown in image below, it is deployed and working fine.



```
root@ip-172-31-23-148: /home/ubuntu/hpa                                          —    □    ×
root@ip-172-31-23-148:/home/ubuntu/hpa# kubectl top nodes
NAME                                   CPU(cores)   CPU%   MEMORY(bytes)   MEMORY%
ip-10-0-0-235.us-west-2.compute.internal   54m      2%     518Mi           15%
root@ip-172-31-23-148:/home/ubuntu/hpa#
```

Now, using simple yaml files, I deployed Nginx and its service in the cluster. Service is used to expose the nginx pod through NodePort.
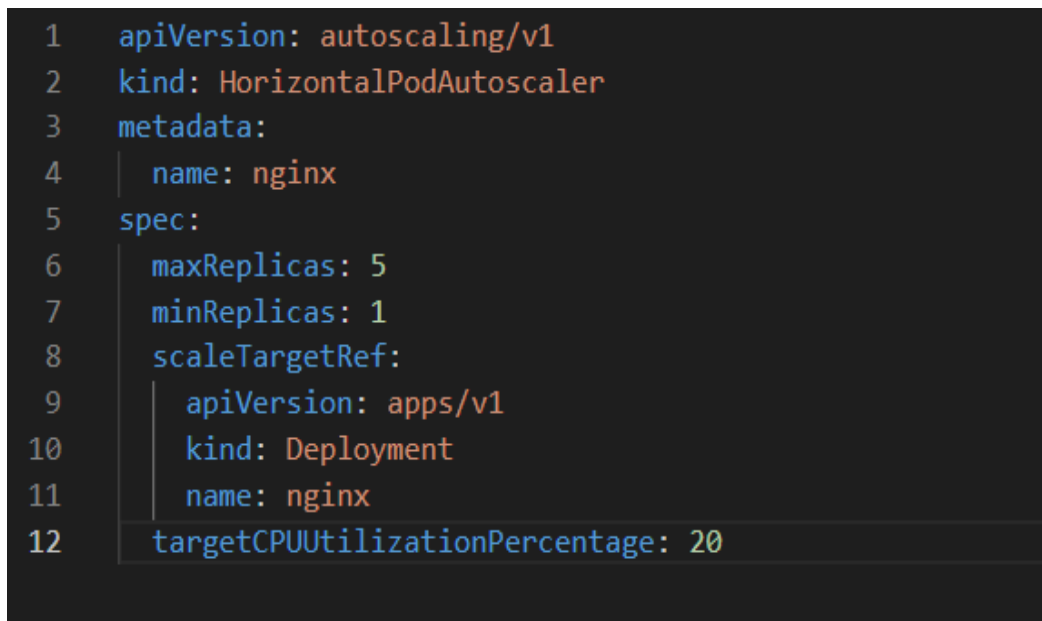
Now, to make stress on the Nginx, I have used the Siege tool for stress testing purpose.

Command to install siege:              apt-get install siege
Command to create stress over Nginx:   siege -q -c 20 -f 2m http://url

For enabling Horizontal Pod Autoscaling, I used the below shown file named hpa.yaml.
Command used: kubectl -f apply hpa.yaml

```
 1   apiVersion: autoscaling/v1
 2   kind: HorizontalPodAutoscaler
 3   metadata:
 4     name: nginx
 5   spec:
 6     maxReplicas: 5
 7     minReplicas: 1
 8     scaleTargetRef:
 9       apiVersion: apps/v1
10       kind: Deployment
11       name: nginx
12     targetCPUUtilizationPercentage: 20
```

As shown in the image below, commands are run and stress on the Nginx is increased.

As it reached the threshold, we set for the horizontal pod autoscaling. It will start creating new pods of the same kind. Hence, our target is accomplished by having more pods in real time to manage the stress on the application.