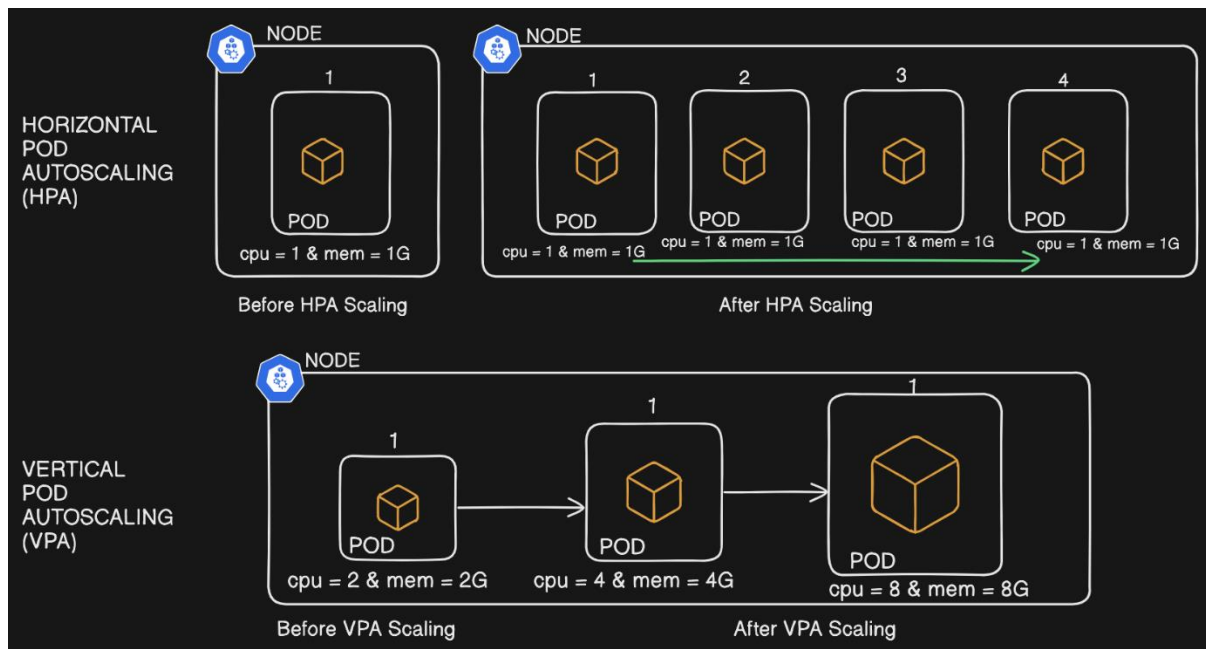


# HORIZONTAL POD AUTOSCALER (HPA) & VERTICAL POD AUTOSCALER (VPA) IN KUBERNETES

**HPA (Horizontal Pod Autoscaler)** and **VPA (Vertical Pod Autoscaler)** are mechanisms that manage the scaling of resources for applications running within the cluster.



## 1. HPA - Horizontal Pod Autoscaler

- **Purpose:** HPA automatically adjusts the number of pod replicas in a deployment, replicaset, or statefulset based on observed CPU, memory usage, or other custom metrics.
- **Usage:** Best suited for stateless applications where multiple instances of the application can run simultaneously without data synchronization issues.
- **When to Use:**
  - ✓ When you expect varying loads and want to handle traffic surges.
  - ✓ For applications where scaling horizontally (adding more pods) is easier than increasing resources of a single pod.
- **Advantages:**
  - ✓ Increased resilience and load balancing since multiple pods can share the workload.
  - ✓ Ideal for applications with distributed or stateless architectures.

## 2. VPA - Vertical Pod Autoscaler

- **Purpose:** VPA automatically adjusts the resource requests and limits (CPU, memory) of existing pods based on observed usage patterns, thereby “scaling up” the resource capabilities of each pod individually.
- **Usage:** Ideal for stateful applications where a single instance with more resources is more efficient than multiple replicas.
- **When to Use:**
  - ✓ For applications where horizontal scaling isn't feasible, such as those requiring shared storage or maintaining a unique state.

# HORIZONTAL POD AUTOSCALER (HPA) & VERTICAL POD AUTOSCALER (VPA) IN KUBERNETES

- ✓ For workloads that require more memory or CPU over time but don't need additional replicas.

## ➤ Advantages:

- ✓ Can help reduce resource underutilization by right-sizing pods.
- ✓ Prevents resource overcommitment by adjusting limits based on actual usage patterns.

## Which One is Better: HPA or VPA?

- **Use Cases:** HPA is generally better for **stateless** and **scalable applications**, while VPA is beneficial for stateful applications or those that benefit from increasing the power of individual instances rather than adding new ones.
- **Combined Approach:** Often, HPA and VPA can be used together. HPA handles the number of replicas, while VPA ensures each pod has the correct resources. However, their combined use should be carefully tested, as scaling horizontally while also changing individual pod sizes can introduce complexity.

## Summary

- Use **HPA** for workloads with fluctuating traffic and load that can be distributed across multiple instances.
- Use **VPA** for applications that benefit from increased resources per pod rather than additional replicas.

## Hands-on (HPA)

### Deployment and Service manifest

```
autoscaling > ! deployyaml
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: php-apache
5  spec:
6    selector:
7      matchLabels:
8        run: php-apache
9    template:
10     metadata:
11       labels:
12         run: php-apache
13     spec:
14       containers:
15         - name: php-apache
16           image: registry.k8s.io/hpa-example
17           ports:
18             - containerPort: 80
19           resources:
20             limits:
21               cpu: 500m
22             requests:
23               cpu: 200m
24     ---
25   apiVersion: v1
26   kind: Service
27   metadata:
28     name: php-apache
29   labels:
30     run: php-apache
31   spec:
32     ports:
33       - port: 80
34     selector:
35       run: php-apache
36
37
38 #to increase the load in CPU run this command, you will get this command in k8's doc --> hpa -- got to bottom --> example --> find increase the load on cpu
39 #kubectl run -i --tty load-generator --rm --image=busybox:1.28 --restart=Never -- /bin/sh -c "while sleep 0.01; do wget -q -O- http://php-apache; done"
40
```

# HORIZONTAL POD AUTOSCALER (HPA) & VERTICAL POD AUTOSCALER (VPA) IN KUBERNETES

Horizontal pod autoscaling manifest, where ever the CPU utilization cross 50% new pod will be created.

```
autoscaling > ! scale.yaml
1  apiVersion: autoscaling/v2
2  kind: HorizontalPodAutoscaler
3  metadata:
4    name: php-apache
5  spec:
6    scaleTargetRef:
7      apiVersion: apps/v1
8      kind: Deployment
9      name: php-apache
10   minReplicas: 1
11   maxReplicas: 10
12   metrics:
13   - type: Resource
14     resource:
15       name: cpu
16       target:
17         type: Utilization
18         averageUtilization: 50
```

```
manoj -->
manoj -->
manoj -->kubectl get pods
No resources found in default namespace.
manoj -->
manoj -->kubectl apply -f deploy.yaml
deployment.apps/php-apache created
service/php-apache created
manoj -->
manoj -->kubectl get pods
NAME                                READY   STATUS    RESTARTS   AGE
php-apache-d87b7ff46-bxk7d         1/1     Running   0           6s
manoj -->
manoj -->kubectl get svc
NAME      TYPE        CLUSTER-IP   EXTERNAL-IP   PORT(S)   AGE
kubernetes ClusterIP  10.96.0.1     <none>        443/TCP      27d
php-apache ClusterIP  10.96.216.60 <none>        80/TCP       17s
manoj -->
manoj -->
```

← pod and service running

HPA Object created

```
manoj -->
manoj -->
manoj -->kubectl autoscale deploy php-apache --cpu-percent=50 --min=1 --max=10
horizontalpodautoscaler.autoscaling/php-apache autoscaled
manoj -->
manoj -->kubectl get hpa
NAME      REFERENCE          TARGETS      MINPODS  MAXPODS  REPLICAS  AGE
php-apache Deployment/php-apache  cpu: <unknown>/50%  1         10        0          11s
manoj -->
manoj -->kubectl get hpa
NAME      REFERENCE          TARGETS      MINPODS  MAXPODS  REPLICAS  AGE
php-apache Deployment/php-apache  cpu: 0%/50%    1         10        1          28s
manoj -->
manoj -->kubectl get hpa
NAME      REFERENCE          TARGETS      MINPODS  MAXPODS  REPLICAS  AGE
php-apache Deployment/php-apache  cpu: 0%/50%    1         10        1          41s
manoj -->
manoj -->
```

← HPA object created for deployment

min of 1 pod, max of 10 pod will be create when load on CPU increases above 50%

after the load, CPU load decreases to below 50% then replicas will downsize to min of 1 pod

# HORIZONTAL POD AUTOSCALER (HPA) & VERTICAL POD AUTOSCALER (VPA) IN KUBERNETES

Applying stress to increase the load on CPU

```
PS C:\Users\manoj_gowda_ac\Desktop\kubernetes>
PS C:\Users\manoj_gowda_ac\Desktop\kubernetes> kubectl run -i --tty load-generator --rm --image=busybox:1.28 --restart=Never -- /bin/sh -c "while sleep
0.01; do wget -q -O- http://php-apache; done"
```

If you don't see a command prompt, try pressing enter.

pod default/load-generator terminated (Error)  
PS C:\Users\manoj\_gowda\_ac\Desktop\kubernetes>

increasing the load on CPU

We can see when CPU utilization crossed 50% new pod got created as we specified in HPA manifest.

```
manoj -->
manoj -->
manoj --># now let me increase the CPU load
manoj -->
manoj -->kubectl get pods
NAME                                READY   STATUS    RESTARTS   AGE
php-apache-d87b7ff46-bxk7d          1/1     Running   0           14m
manoj -->
manoj -->kubectl get hpa
NAME      REFERENCE                TARGETS      MINPODS   MAXPODS   REPLICAS   AGE
php-apache  Deployment/php-apache     cpu: 0%/50%   1         10        1           7m12s
manoj -->
manoj -->kubectl get hpa --watch
NAME      REFERENCE                TARGETS      MINPODS   MAXPODS   REPLICAS   AGE
php-apache  Deployment/php-apache     cpu: 79%/50%   1         10        1           7m18s
php-apache  Deployment/php-apache     cpu: 250%/50%  1         10        2           7m30s
php-apache  Deployment/php-apache     cpu: 138%/50%  1         10        4           7m46s
```

when the load on CPU increase above 50%, we can see the HPA got activated

```
manoj -->kubectl get pods
NAME                                READY   STATUS    RESTARTS   AGE
load-generator                      1/1     Running   0           56s
php-apache-d87b7ff46-86v88          1/1     Running   0           22s
php-apache-d87b7ff46-bxk7d          1/1     Running   0           15m
php-apache-d87b7ff46-dfk2f          1/1     Running   0           6s
php-apache-d87b7ff46-rzb4m          1/1     Running   0           22s
php-apache-d87b7ff46-tvb6n          1/1     Running   0           37s
manoj -->
manoj -->
```

increased pod because of increase in load on CPU which cross above 50%