# An academic look at
# KUBERNETES AUTO-SCALING

State-of-the-art
Approaches
Challenges
Innovation

## HORIZONTAL POD AUTO-SCALING (HPA)

The most basic & common form of auto-scaling provided natively by Kubernetes.
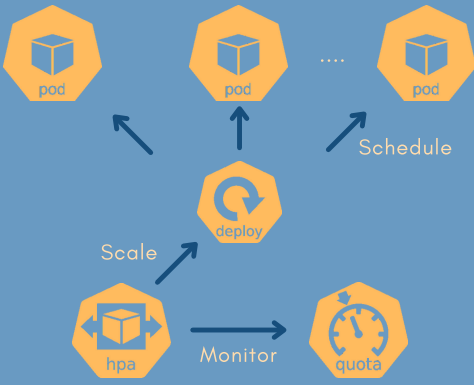
Paper: *Agnostic Approach for Microservices Autoscaling in Cloud Applications*

Author: Khaleq et al.

**1**

## K8S HPA ≃ HORIZONTAL SCALING

Kubernetes HPA varies slightly from the generic Horizontal scaling concept known outside its ecosystem. Instead of varying the number of nodes to scale in or out, applications in Kubernetes are scaled by scheduling or removing Pods, which are the building block of our systems and the unit containing our virtualized instances.

**2**



**3**

In its most basic form, all HPA does is periodically monitor resources (CPU, Memory) quotas configured. It will then call the controller to scale in or out, matching the observed metrics and following the specified scaling behaviour.

```
1 type: Resource
2 resource:
3   name: cpu
4   target:
5     type: Utilization
6     averageUtilization: 60
```

```
1 behavior:
2   scaleDown:
3     policies:
4       - type: Pods
5         value: 4
6         periodSeconds: 60
7       - type: Percent
8         value: 10
9         periodSeconds: 60
```

**4**

## CUSTOM METRICS

Custom Metrics are a complimentary feature to HPA and other auto-scaling techniques.

Paper: Horizontal Pod Autoscaling in Kubernetes for Elastic Container Orchestration

Author: Nguyen et al.

**5**

Custom Metrics are widely used nowadays. It provides for a lot of flexibility in the configuration of auto-scaling. Pretty much anything can be fed to the Kubernetes API as a metric. The API allows us to define & declare metrics and consume data from external & internal sources, such as the Kubernetes Metrics Server.

```
1 type: Object
2 object:
3   metric:
4     name: http_requests
5     selector: {matchLabels:
   {verb: GET}}
```

```
1 - type: External
2   external:
3     metric:
4       name: queue_messages_ready
5       selector:
6         matchLabels:
7           queue: "worker_tasks"
8     target:
9       type: AverageValue
10      averageValue: 30
```

**6**

By using adapters, we can easily transfer data from Prometheus, for example, to any federated monitoring solution such as Stack driver and by that improving accessibility and metrics management.

```
1 spec:
2   containers:
3     - name: app
4     ....
5     - name: prometheus-to-sd
6       image: gcr.io/google-containers/prometheus-
   to-sd:v0.5.0
7       command: ["/monitor"]
8       args:
9         - --source=http://localhost:8080
10        - --stackdriver-prefix=custom.googleapis.com
11        - --pod-id=$(POD_ID)
12        - --namespace-id=$(POD_NAMESPACE)
```
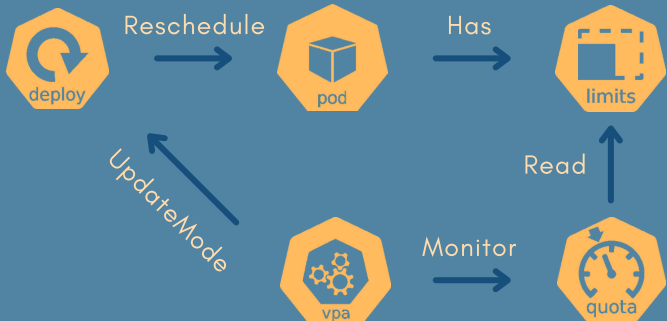
**7**

## VERTICAL POD AUTO-SCALING (VPA)

Kubernetes VPA is far from the generic notion of Vertical Scaling. The feature itself is still in development. Also, due to its nature, VPA is a platform specific feature, meaning that you might not be able to use it depending on your cloud provider or environment.

Paper: Exploring Potential for Non-Disruptive Vertical Auto Scaling and Resource Estimation in Kubernetes

Author: Rattihalli et al.

**8**



**9**

## K8S VPA ≠ VERTICAL SCALING

VPA in Kubernetes consists of editing resources requests and limits of Pods on the fly and therefore removes the need to specify them manually. It keeps track of resource usage, and then decides to increase or decrease the configured values for newly scheduled Pods, or for existing ones in some cases. The update process still requires to reschedule the Pod for now, but it is planned to have an *UpdateMode* without this constraint implemented at some point.

**10**

VPA configuration is still pretty similar to HPA. Custom Metrics are used in the same way, and instead of specifying a replicas number bounds, we specify resource bounds and target.

```
1 apiVersion: autoscaling.k8s.io/v1
2 kind: VerticalPodAutoscaler
3 metadata:
4   name: my-vpa
5 spec:
6   targetRef:
7     apiVersion: "apps/v1"
8     kind:       Deployment
9     name:       my-rec-deployment
10  updatePolicy:
11    updateMode: "Off"
```

```
2 recommendation:
3   containerRecommendations:
4     - containerName: my-container
5       lowerBound:
6         cpu: 536m
7         memory: 262144k
8       target:
9         cpu: 587m
10        memory: 262144k
11      upperBound:
12        cpu: 27854m
13        memory: "545693548"
```

**11**

## CLUSTER AUTO-SCALING

Cluster auto-scaling is closer to the conventional Horizontal Scaling known outside K8s ecosystem. Being directly bound to your system's infrastructure, it is a platform specific feature. In fact, in most cases, this would be configured directly through your cloud provider interface.

Paper: An Experimental Evaluation of the Kubernetes Cluster Autoscaler in the Cloud

Author: Tamiru et al.

**12**

As you may have guessed, Cluster auto-scaling allows for the dynamic addition and removal of nodes from a Kubernetes cluster. It also offers additional configurations, such as providing differentiated node pools and a selection strategy to fit into more specific use-cases. It also comes with a profile configuration option, allowing for cost optimization, for example.

```
$ gcloud container clusters create example \
    --num-nodes 2 \
    --zone us-central1-a \
    --enable-autoscaling --min-nodes 1 --max-nodes 4
```

```
$ gcloud beta container clusters update example \
  --autoscaling-profile optimize-utilization
```

**13**

## MULTIDIMENSIONAL POD AUTO-SCALING (MPA)

Now this is where we start truly getting out of the beaten tracks. MPA is still super early in its development stage. Even finding documentation for it might be a hustle in some cases. At this point, I would not consider it of much use with all the limitations it has. It is, however, a very interesting concept that might have a lot of potential later in its road map.

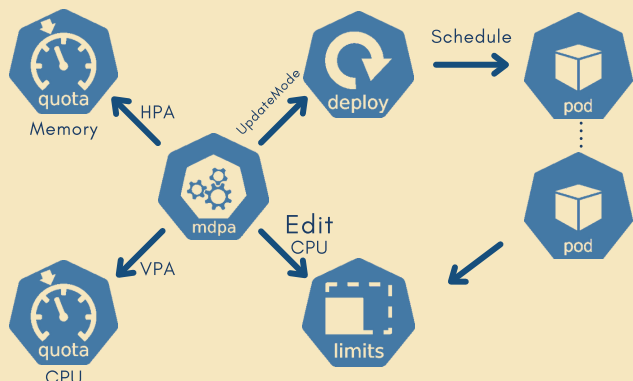Paper: Adaptive scaling of Kubernetes pods

Author: Balla et al.

**14**

## K8S MPA = VPA || HPA

The idea behind MPA isn't that complicated. It is basically an auto-scaling technique, allowing for both horizontal AND vertical auto-scaling simultaneously. However, HPA being coupled to Pod limits to trigger scaling and VPA editing these same limits, encountering conflicts and scaling misbehaviour, is inevitable. MPA dodges the issue by limiting the type of resource each dimension gets to work with.

```
1 apiVersion: autoscaling.gke.io/v1beta1
2 kind: MultidimPodAutoscaler
3 metadata:
4   name: autoscaler
5 spec:
6   scaleTargetRef:
7     apiVersion: apps/v1
8     kind: Deployment
9     name: app
```

**15**



**16**

For now, MPA manages CPU though its vertical axis and memory through its horizontal. It is however planned to have an implementation solving this limitation soon enough.

```
2 goals:
3   metrics:
4     - type: Resource
5       resource:
6         name: cpu
7         target:
8           type: Utilization
9           averageUtilization: 60
```

```
2 constraints:
3   global:
4     minReplicas: 1
5     maxReplicas: 5
6     containerControlledResources: [ memory ]
7     container:
8       - name: '*'
9         requests:
10          minAllowed:
11            memory: 1Gi
12          maxAllowed:
13            memory: 2Gi
14      policy:
15        updateMode: Auto
```

**17**

MED MOUINE
E: MOHAMED.MOUINE.2@ULAVAL.CA