

134 Vertical Pod Auto-scaling (VPA) vertical auto scaling

Scaling resources for a workload the manual way

Manually scaling workload vertically.

As an sys administrator tasked to make sure that there is always sufficient workload to support demand for this application.

```
nginx.yaml
```

apiVersion: apps/v1
kind: Deployment
metadata:
 name: my-app
spec:
 replicas: 1
 selector:
 matchLabels:
 app: my-app
 template:
 metadata:
 labels:
 app: my-app
 spec:
 containers:
 - name: my-app
 image: nginx
 resources:
 requests:
 cpu: "250m"
 limits:
 cpu: "500m"

And from the deployment configuration, I see that this pod requests 250 millicores of CPU & has a limit of 500 millicores of CPU, this means that 500 millicores is the max capacity it gets

<code>kubectl top pod my-app-pod</code>	NAME	CPU (cores)	MEMORY (bytes)
	my-app-pod	450m	350 MiB

we would run '`kubectl top pod`' command & monitor the resource consumption of the pod

Remember we need to have metrics server running on the cluster to be able to monitor, the same as before & when it reaches a specific threshold what I'm going to do to scale the pod vertically is to run the `kubectl edit deployment` command, we then run the `kubectl edit deployment` command & change the resource requests & limits assigned to the deployment under the pod template, under the `containers` section & then save it

`kubectl edit deployment my-app`

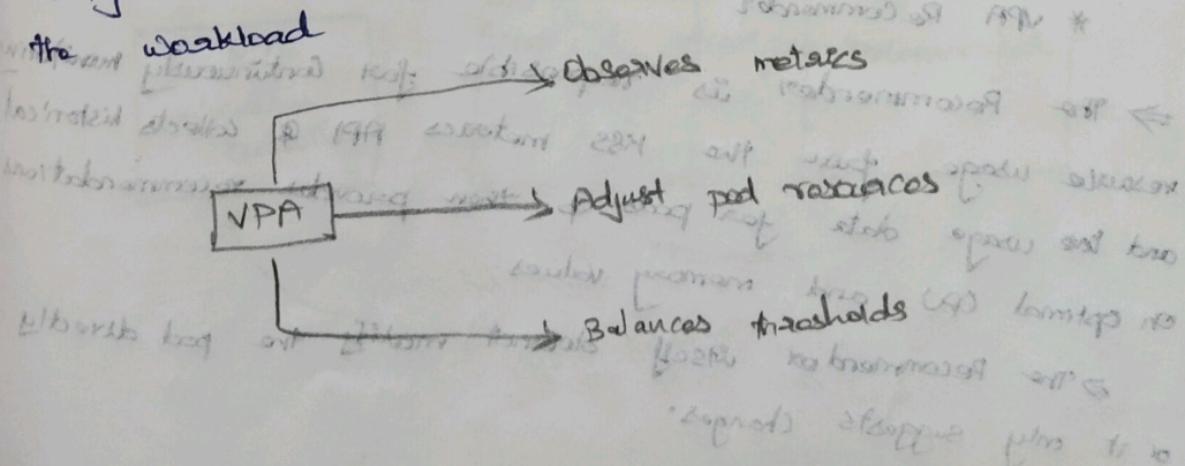
what would happen?

If it would kill that pod & create a new pod

This is manual approach.

For automation we have Vertical Pod Autoscaler

similar to the Horizontal Pod Autoscaler, the Vertical Pod Autoscaler will continuously monitors the metrics & then it automatically increases or decreases the resources assigned to the pods in a deployment & thus balances the workload



Look at it instead

So unlike the Horizontal Pod Autoscaler, the Vertical Pod Autoscaler do not come built-in.

As such, we must deploy it, so, we first apply the Vertical Pod AutoScaler definition file available in the GitHub repo.

Then apply -f <https://github.com/k8s-autoscaler/releases/latest/> download) vertical-pod-autoscaler.yaml

And then

vpa-admission-controller-xxx	Running
vpa-recommender-xxxx	Running
vpa-updater-xxxx	Running

Search for vpa, we will be able to see that there are multiple components deployed, so there is the admission-controller recommender and the updater service which should be running

lets look at what these are

so, that VPA deployment consists of multiple components set up as follows:

- * VPA admission controller
- * VPA updater
- * VPA Recommender

⇒ The Recommender is responsible for continuously monitoring resource usage from the K8s metrics API & collects historical and live usage data for pods, & then provides recommendations on optimal CPU and memory values

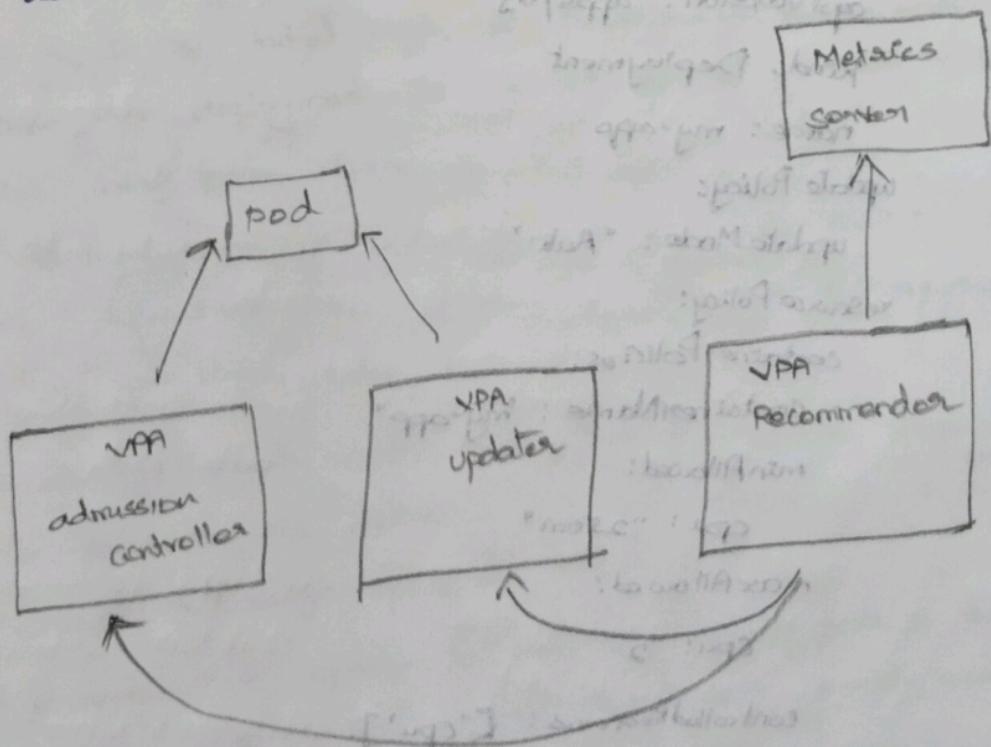
⇒ The Recommender itself does not modify the pod directly or it only suggests changes.

The Updater detects pods that are running with suboptimal resources & evicts them when an update is needed.

It gets the information from the Recommender & monitors the pod & if a pod needs to be updated, it evicts them. So that means it just terminates the pod and the admission controller involves the pod creation process & uses the recommendations from the Recommender again to then mutate the pod spec to apply the recommended CPU & memory values at startup.

This ensures that the newly created pods start with the correct resource requests

Basically the VPA Recommender collects information, the updater monitors or gets the information from the Recommender, compares that to the actual pod, & if the pod is beyond the threshold, it kills the pod. and whether it kills or not that depends on the policy that we will talk about in a bit.



But, ideally it would kill the pod, & then the Admission Controller intervenes because when a pod is killed, it's automatically...
when the deployment will automatically recreate the pod.
And when it does that, the Admission Controller intervenes
& updates the resources so that the pod now comes up
with a new size.

Create a VPA Service using deployment file

Note: There is no imperative command to create a VPA
because it's not a built-in component like the HPA

my-app-vpa.yaml

apiVersion: autoscaling.k8s.io/v1

kind: VerticalPodAutoscaler

metadata:

name: my-app-vpa

spec:

targetRef:

apiVersion: apps/v1

kind: Deployment

name: my-app

updatePolicy:

updateMode: "Auto"

resourcePolicy:

containerPolicies:

- containerName: "my-app"

minAllowed:

cpu: "250m"

maxAllowed:

cpu: "2"

controlledResources: ["cpu"]

There are 4 modes that VPA operates in

Mode

Description

off

Only recommends. Does not change anything

Initial

Only changes on Pod Creation. Not later

Recreate

Evicts pods if usage goes beyond range

Auto

Updates existing pods to recommended numbers.
For now this behaves similar to Recreate. But
when support for "Inplace update of Pod
Resources" is available that mode will be
preferred

⇒ The off mode, which only recommends changes but
does not do anything. In this case, it's only the recommender
working. The "updater" and the admission controller are not
working in case of the off mode.

⇒ The initial mode, that only changes pods on creation,
not later

⇒ So in the initial mode, the Recommender recommends a change
& when the deployment is scaled for some other reason
and the new pods come up.

⇒ The admission controller intervenes & changes the pod's
resource definitions. The update in this case does not work
to kill as, you know, take down the pod.

⇒ In case of Recreate mode, that's where the updater component
jumps in & evicts an existing pod, when its resource consumption
goes beyond a range specified

⇒ Auto mode, this updates to existing ~~number~~ pods to the
recommended numbers. Again, for now, auto mode works exactly
as the Recreate mode. But when the support for the in-place
update of pod resources is available, the one we discussed

In previous lecture then that mode would be preferred

Over recreate. So, "automate" is built for future

→ check if the update of pod resources are available in stable version of k8s, that would be the default behavior. But as of now, the Automate works exactly like Recreate, when a pod goes beyond a range, it is killed.

→ Now, VPA will monitor resource usage & suggest adjustments, so we can check the recommendations with the

Kubectl describe vpa my-app-vpa

Recommendations:

Target:

CPU: 1.5

but currently consumer also thinks you can see the recommendation. So VPA suggests increasing CPU to 1.5

when to use VPA / HPA

Features	VPA	HPA
Vertical Scaling	Increases CPU & memory of existing pods	Adds/removes pods based on loads
Method	existing pods are scaled up or down	new pods are added or removed
Pod Behavior	restarts pods to apply new resource values	keeps existing pods running
Handles Traffic spikes	No, because scaling requires a pod restart	yes, instantly adds more pods
optimizes costs	prevents over-provisioning of CPU/MEMORY	Avoids unnecessary idle pods
Best for	stateful workloads, CPU/memory-heavy apps (DBs, ML workloads)	webapps, microservices, stateless services
Example	Databases (MySQL, PostgreSQL)	web services (Nginx, API services)
Workload	JVM based apps, AI/ML workloads	message queues, microservices

Summary

- ⇒ VPA is about optimizing resource allocation for individual pods
- ⇒ HPA is about scaling many pods dynamically based on demand

⇒ Hence choosing the right auto-scale depends on the workload type & how you need to scale your apps.

VPA quick notes

what is VPA

VPA automatically adjusts CPU and memory requests/limits for pods based on usage over time

- * Helps avoid under-provisioning (OOM/crash) & over-provisioning (wasted resources)
- * VPA adjusts pod size, not the no. of pods.

what VPA adjusts

Resource Field
resources.requests

resources.limits

Replica Count

Adjusted by VPA

yes

Notes
key input for scheduling

sometimes optional - may be configured.

NO option we use HPA for this

VPA modes

Mode

OFF

only shows recommendations; no change to running pods

Initial

Applies resources only when pod is created
Initially

Auto

updates pods automatically by evicting & restarting them

Internal VPA Components

Component	Role
Recommender	Watches usage, calculates resource suggestions
Updater	Evicts pods if they fall outside recommended ranges
Admission Controllers	Injects recommended values into new pods

Typical VPA workflow

- 1) App runs with set requests / limit
- 2) Recommender observes resource usage
- 3) Recommendations appear in kubelet describe vpa
- 4) Updater (in Auto) evicts pods to apply new settings
- 5) New pods start with updated resources

VPA YAML Example

```
apiVersion: autoscaling.k8s.io/v1
kind: VerticalPodAutoscaler
metadata:
  name: myapp-vpa
spec:
  targetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: myapp
  updatePolicy:
    updateMode: Auto # or Initial, off
```

Watch Out for these

Common Problems + Fixes	Fix
symptom No recommendations	Likely Cause matrices-server not running Install matrices-server
pods restart unexpectedly	VPA in Auto mode switch user to initial (or) off
HPA + VPA conflict	Both control CPU requests use HPA for CPU & VPA for memory
Recommendations too high	No resource limits use maxAllowed field in VPA
short-lived pods, no data	Not enough runtime use long-running apps
NPA vs HPA	vs.
Feature	VPA
Adjusts CPU	yes
Adjusts memory	yes
Restarts pods	yes
Adjusts pod count	No Yes
	Both Together only one should handle it safe to combine N/A Combine with care

- Safe VPA + HPA
- ⇒ Use **HPA** for scaling **replicas** (e.g. based on CPU)
 - ⇒ Use **VPA** for memory tuning
 - ⇒ Prevents conflicts

Resource Policy:

Container Policies:

- containerName: my-app

controlled Resources: ["memory"] # not CPU

Key Commands

Task

View VPA Status

kubectl describe vpa <name>

List all VPAs

kubectl get vpa

Create VPA from yaml

kubectl apply -f vpa.yaml

View VPA component pods

kubectl get pods in kube-system

Final Summary

- ⇒ VPA makes pods smarter by tuning their CPU/Memory
- ⇒ It learns from metrics server
- ⇒ It can restart pods (Auto) or just suggest off
- ⇒ Use with HPA carefully - avoid overlapping

control of resource.