

Daniele Polencic — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
@danielepolencic



When designing a Kubernetes cluster, you might need to answer questions such as:

How long does it take for the cluster to scale?

How long do I have to wait before a new Pod is created?

Let's explore...

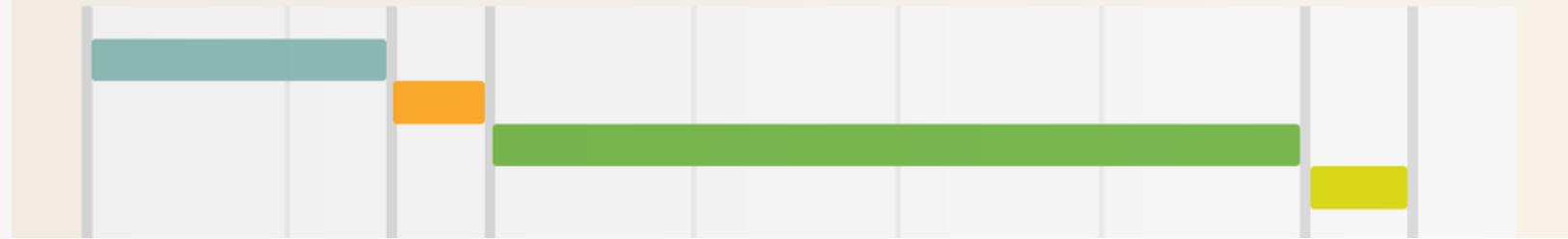
# COMBINING AUTOSCALERS in KUBERNETES

1m30s

30s

4m

30s



9:00 PM · Sep 19, 2023

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic



1/

Four significant factors affect scaling:

- ① Horizontal Pod Autoscaler reaction time
- ② Cluster Autoscaler reaction time
- ③ Node provisioning time
- ④ Pod creation time

#### FACTOR AFFECTING SCALING

## 1 **Horizontal Pod Autoscaler**

The time necessary to detect a scaling opportunity

## 2 **Cluster Autoscaler**

The reaction time for a scaling increment

## 3 **Node provisioning**

The time necessary to create a worker node

## 4 **Pod creation**

The time necessary to download a container image and start the pod.

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic

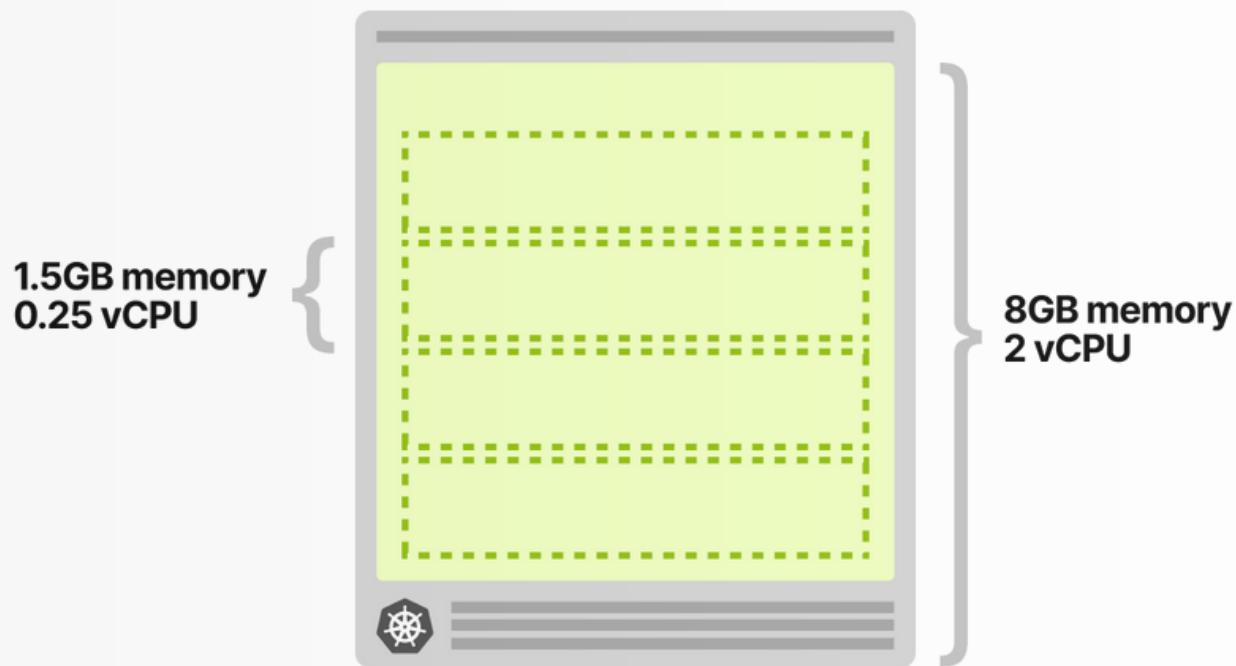


2/

Let's have a look at an example

Imagine having an application that requires and uses 1.5GB of memory and 0.25 vCPU at all times

You provisioned a cluster with a single node of 8GB and 2 vCPU, and 4 replicas; what happens next?



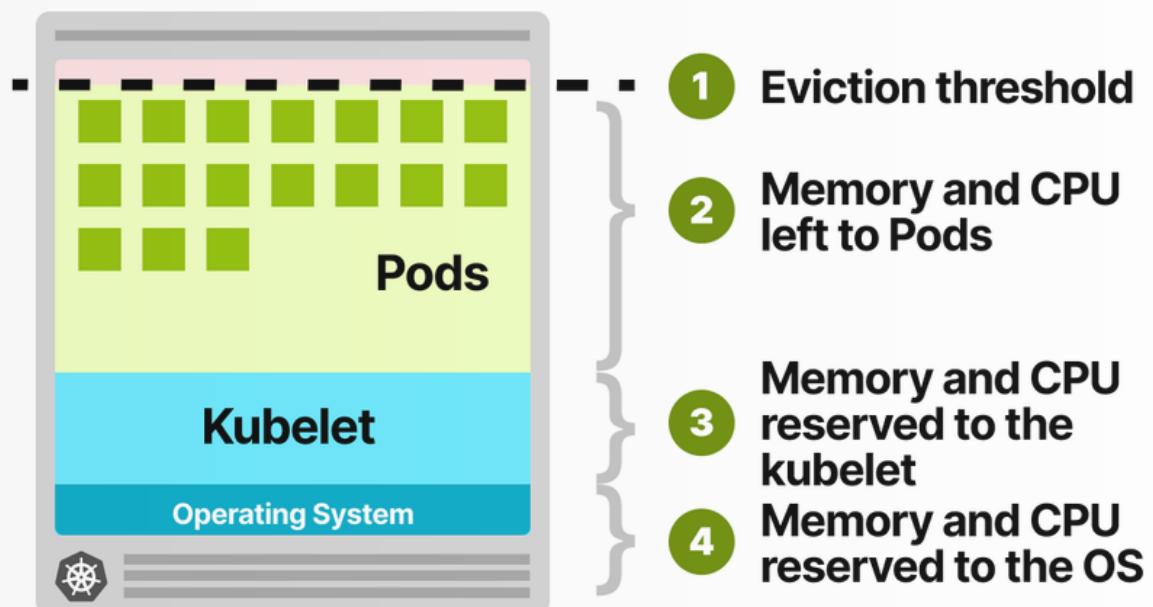
Daniele Polencic — [@danielepolencic](mailto:@danielepolencic@hachyderm.io)



3/

The pod is pending because you still need to consider reserved resources

Not all available memory and CPU are usable by pods



9:02 PM · Sep 19, 2023

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic

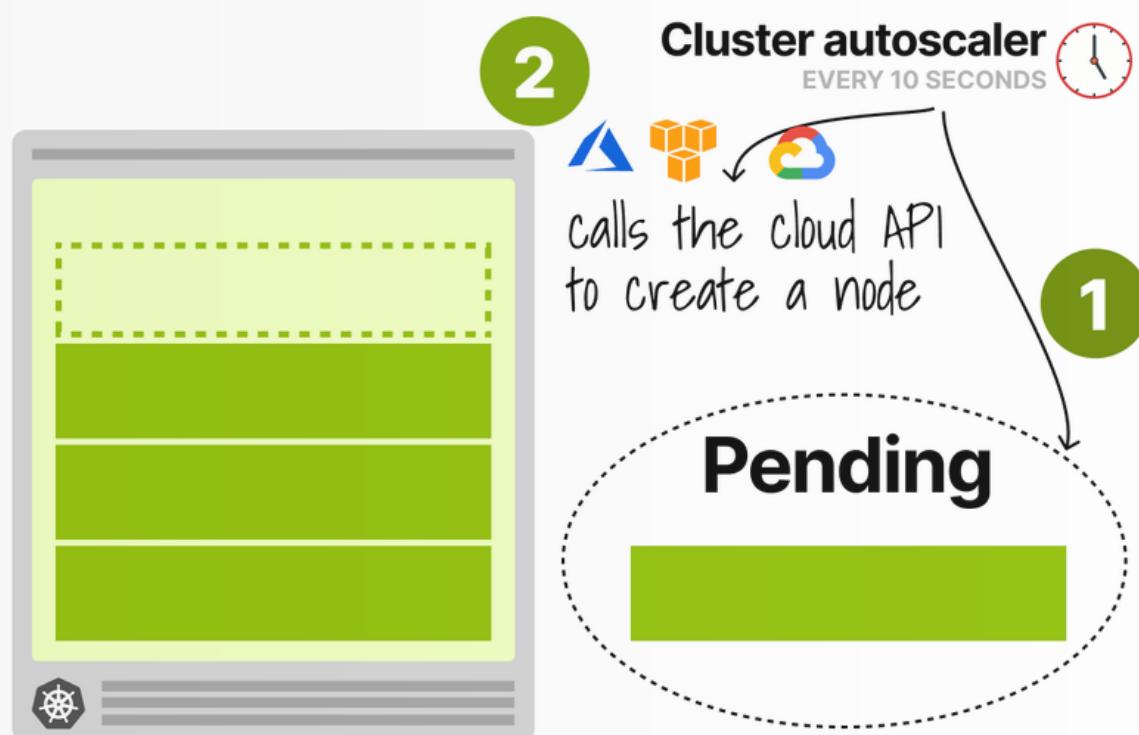


4/

But why does the cluster autoscaler wait for the Pod to be pending before it triggers creating a node?

The CA doesn't measure CPU and memory

Instead, it checks for unschedulable Pods every 10 seconds; if there's at least one, it creates a node



Daniele Polencic — @danielepolencic@hachyderm.io

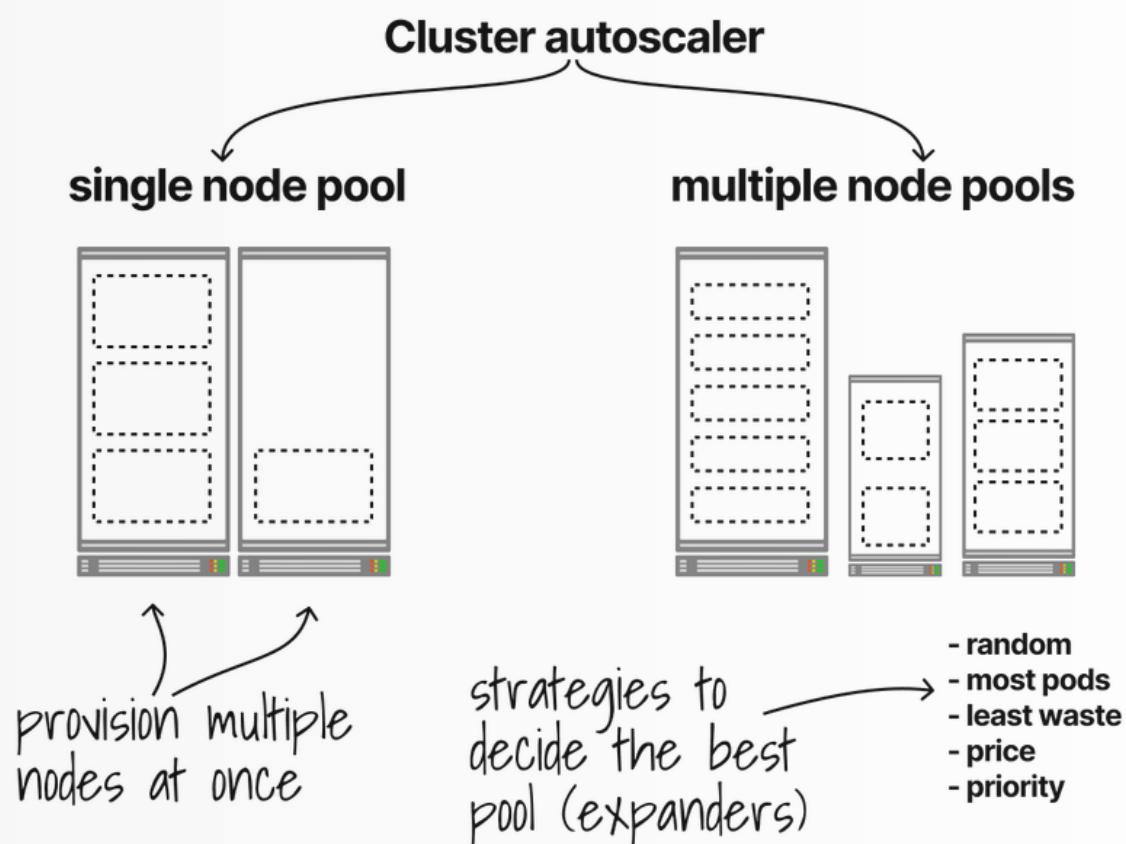
@danielepolencic



5/

It can create more than one node simultaneously if you have several pods

If you have several node pools, you can choose a strategy (expander) to decide which pool is selected

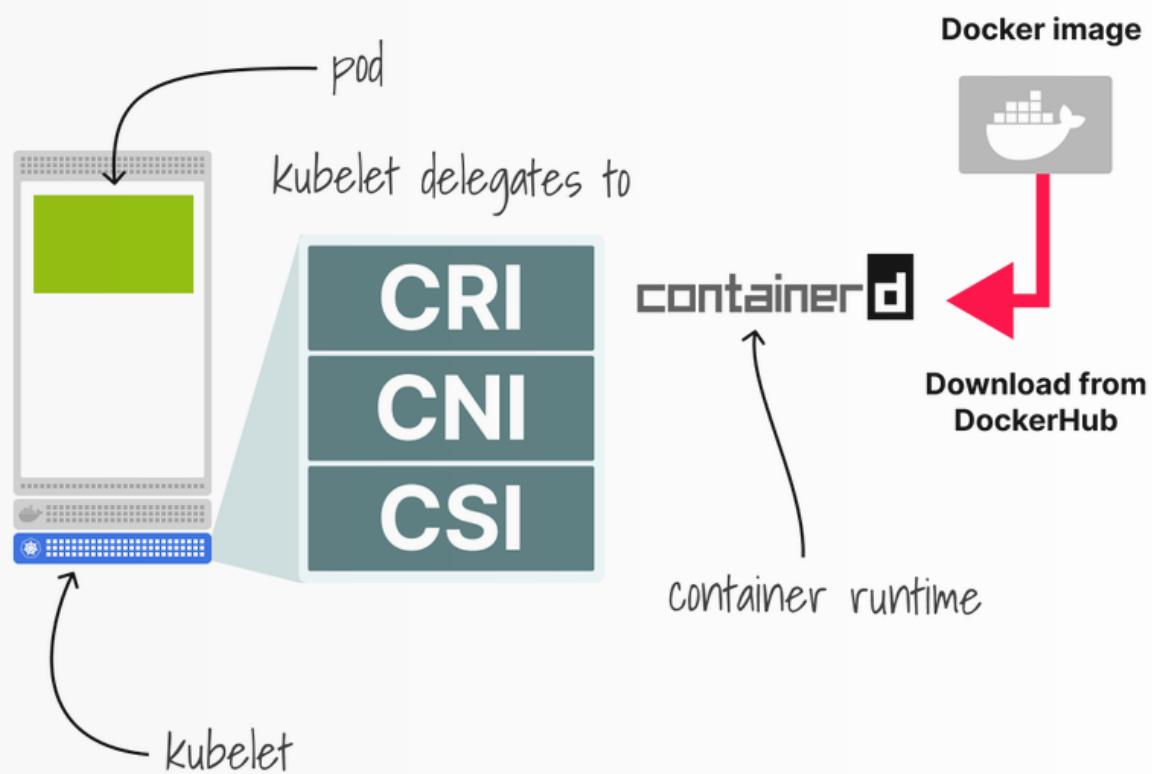


Daniele Polencic — [@danielepolencic](mailto:@danielepolencic@hachyderm.io)



6/

As a final step, the node is provisioned, the container image is pulled from the registry and the pod is launched



9:04 PM · Sep 19, 2023

Daniele Polencic — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
@danielepolencic



7/

Now that you're familiar with the process let's look at how this translates to timings

The Horizontal Pod Autoscaler might take up to 1m30s to increase the number of replicas

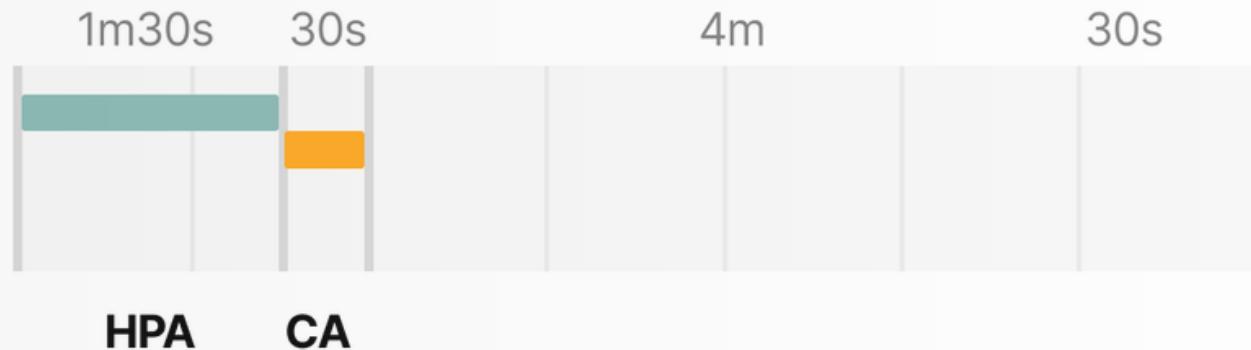


Daniele Polencic — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
@danielepolencic



8/

The Cluster Autoscaler should take less than 30 seconds for a cluster with less than 100 nodes and less than a minute for a cluster with more than 100 nodes



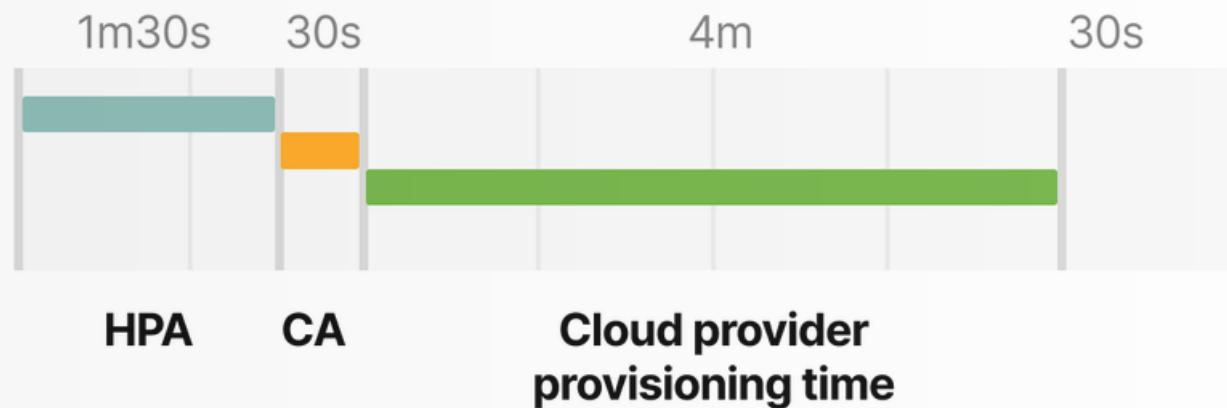
9:04 PM · Sep 19, 2023

Daniele Polencic — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
@danielepolencic



9/

The cloud provider might take 3 to 5 minutes to create the computer resource



9:05 PM · Sep 19, 2023

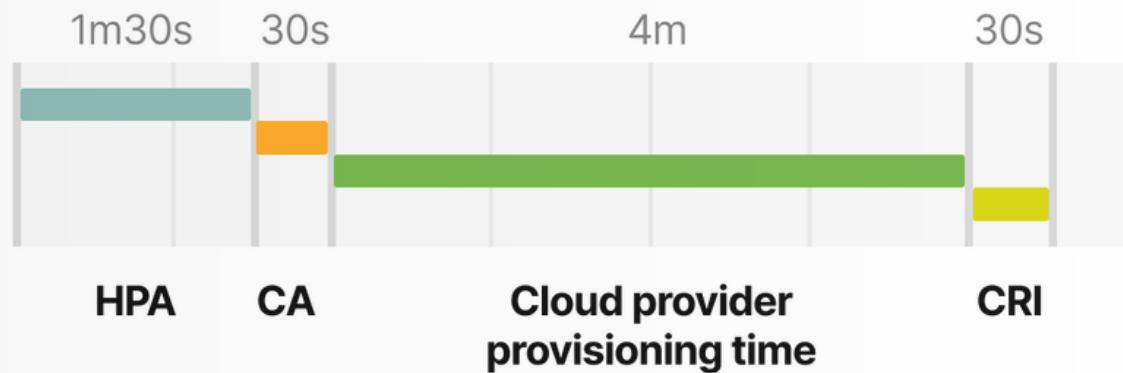


Daniele Polencic — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
[@danielepolencic](https://twitter.com/danielepolencic)



10/

The container runtime could take up to 30 seconds to download the container image



9:05 PM · Sep 19, 2023

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic



11/

In the worst case, with a small cluster, it takes 6m30s from when a pod is pending to launch it in a new node

If the cluster has more than 100 nodes, it's up to 7m

<b>HPA delay</b>	<b>1m30s</b>	worst case lead time for scaling pods with a cluster <100 nodes
<b>CA delay</b>	<b>30s</b>	
<b>Cloud provider delay</b>	<b>4m</b>	
<b>CRI</b>	<b>30s</b>	
<b>Total</b>		<b>6m30s</b>

<b>HPA delay</b>	<b>1m30s</b>	worst case lead time for scaling pods with a cluster 100+ nodes
<b>CA delay</b>	<b>1m</b>	
<b>Cloud provider delay</b>	<b>4m</b>	
<b>CRI</b>	<b>30s</b>	
<b>Total</b>		<b>7m</b>

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic



12/

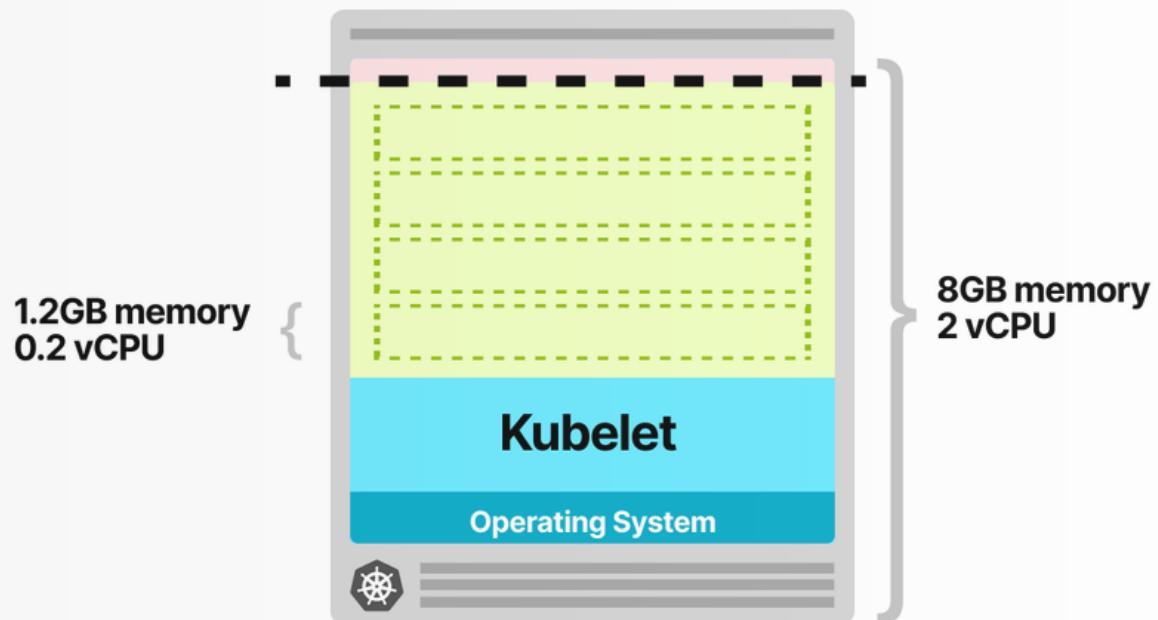
How can you fix this?

If you can, don't scale

Select a node instance type that is optimized for your workload

You can use this calculator to help you with the task:

[learnk8s.io/kubernetes-ins...](https://learnk8s.io/kubernetes-instance-calculator/)



Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic



13/

If your traffic follows predictable patterns, scale your fleet up and down with the CronHPA autoscaler [github.com/AlibabaCloudContainerService/keda.sh](https://github.com/AlibabaCloudContainerService/keda)

With those, you can control the number of replicas with a cron expression



```
apiVersion: autoscaling.alibabacloud.com/v1beta1
kind: CronHorizontalPodAutoscaler
metadata:
  labels:
    controller-tools.k8s.io: "1.0"
  name: cronhpa-sample
spec:
  scaleTargetRef:
    apiVersion: apps/v1beta2
    kind: Deployment
    name: my-deployment
  jobs:
    - name: "scale-down"
      schedule: "@date 2023-09-28 23:59:00"
      targetSize: 1
    - name: "scale-up"
      schedule: "@date 2023-09-30 00:00:00"
      targetSize: 3
```

scale down

scale up

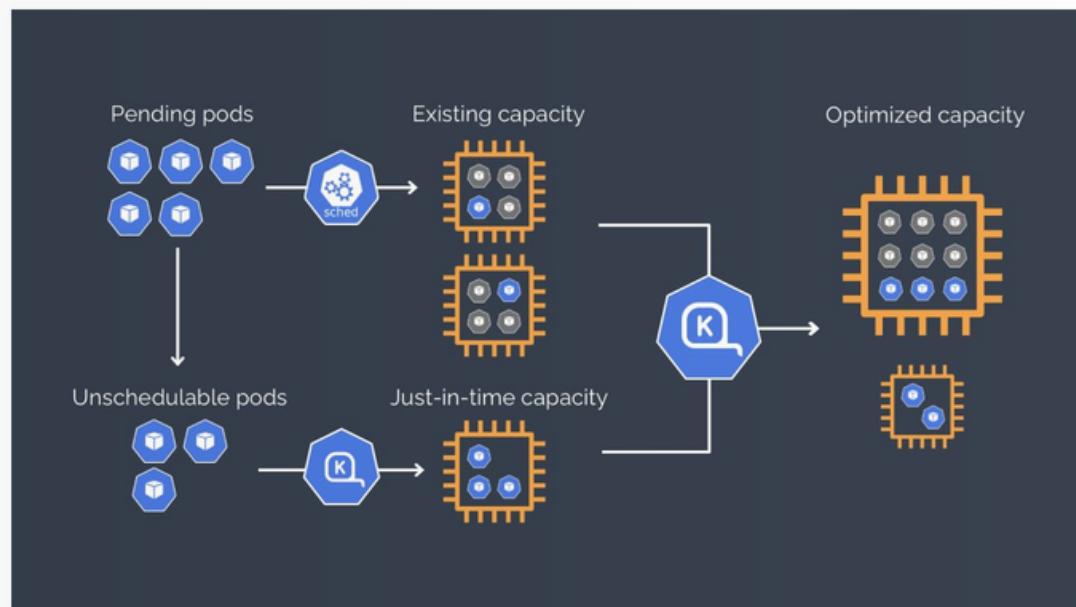
Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic



14/

If you are on AWS, use Karpenter [github.com/aws/karpenter](https://github.com/aws/karpenter)

Karpenter uses a more proactive model that can provision the correct instance and minimize provisioning lead time



9:06 PM · Sep 19, 2023

Daniele Polencic — @danielepolencic@hachyderm.io  
@danielepolencic

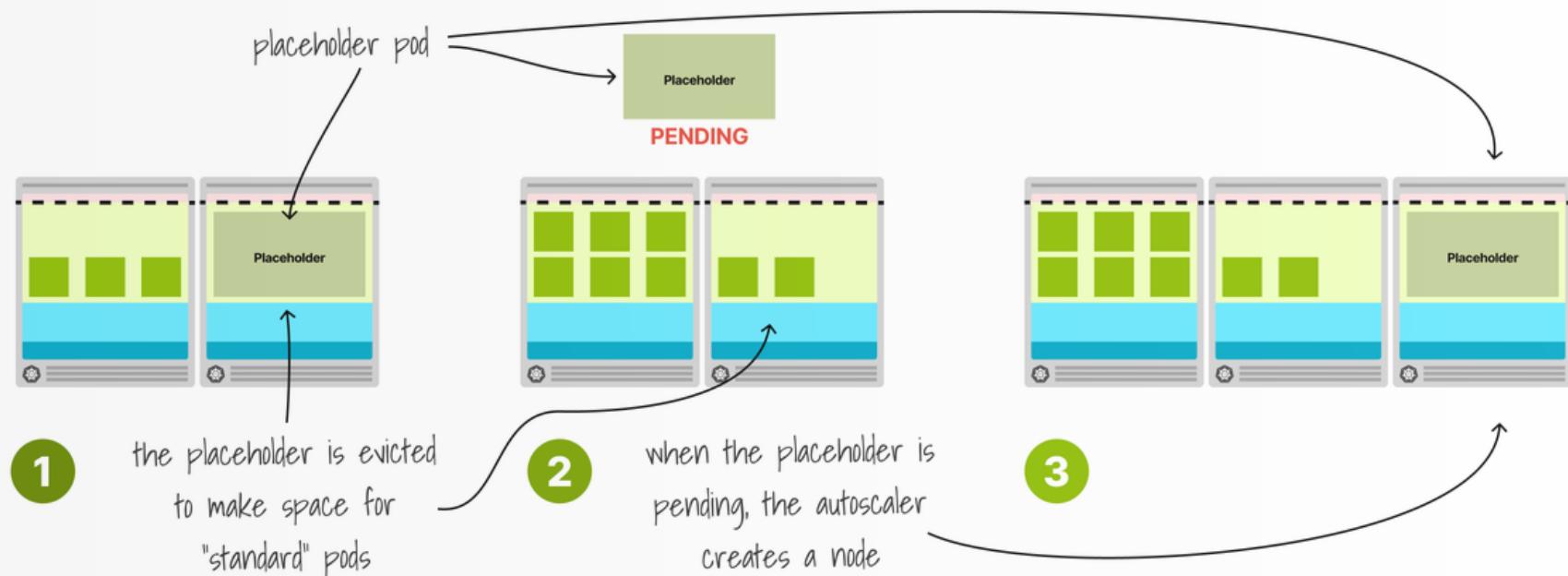


15/

If everything else fails, you can use a placeholder pod always to keep a spare node in the cluster

The pod is configured with the lowest priority and evicted when another pod needs the space

[github.com/kubernetes/aut...](https://github.com/kubernetes/aut...)





Daniele Polencic — [@danielepolencic](mailto:@danielepolencic@hachyderm.io)



16/

This thread is based on [@SoulmanIqbal](#)'s webinar on combining autoscalers

He will present it live next Thursday as the second episode of "Cost optimization and efficiency in Kubernetes" organized by [@Akamai](#) and [@learnk8s](#)

Sign up here (it's free): [event.on24.com/eventRegistration](http://event.on24.com/eventRegistration)...

# KUBERNETES (DOWN)SCALING COMBINING AUTOSCALERS FOR MINIMAL RESOURCE ALLOCATIONS

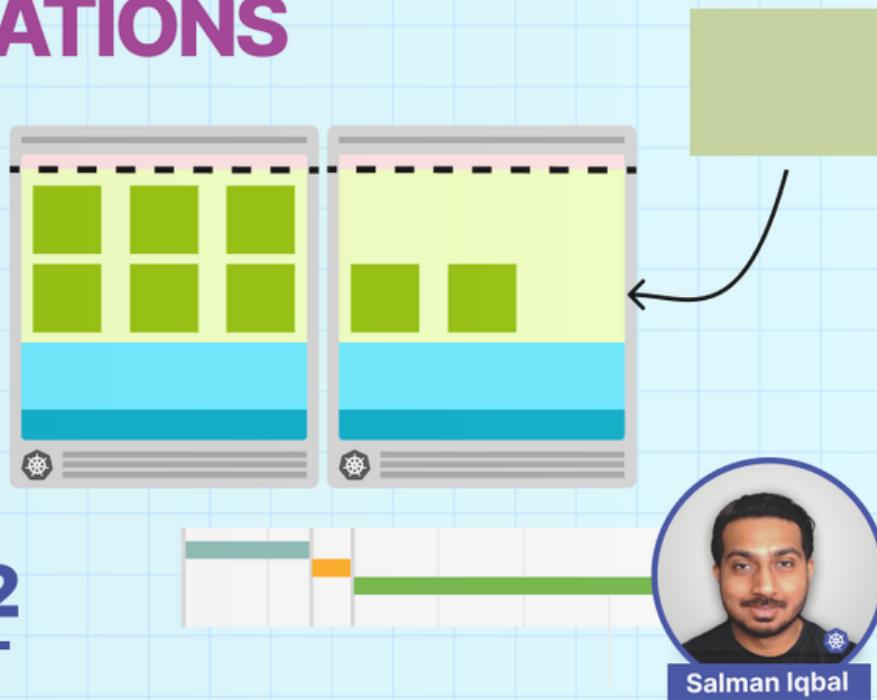
28th of Sep

8am PT | 5pm CET



REGISTER HERE

[bit.ly/k8s-optimize-2](http://bit.ly/k8s-optimize-2)



Salman Iqbal

9:09 PM · Sep 19, 2023



**Daniele Polencic** — [@danielepolencic@hachyderm.io](mailto:@danielepolencic@hachyderm.io)  
[@danielepolencic](https://twitter.com/danielepolencic)



17/

And finally, if you've enjoyed this thread, you might also like:

- The Kubernetes workshops that we run at Learnk8s [learnk8s.io/training](https://learnk8s.io/training)
- This collection of past threads [twitter.com/danielepolencic...](https://twitter.com/danielepolencic/)
- The Kubernetes newsletter I publish every week [learnk8s.io/learn-kubernetes-newsletter...](https://learnk8s.io/learn-kubernetes-newsletter)

9:10 PM · Sep 19, 2023 · 646 Views