

76. Priority classes

101

Opt 2.07

⇒ Kubernetes runs different apps as pods with different priorities. For eg:- the Kubernetes Control Plane Components (apiserver, scheduler, etc.) itself run as pods within the cluster & no matter what those workloads need to always run.

⇒ So, we may also run databases, critical apps, running on the cluster as well as lower priority workloads like background jobs.

⇒ Now, we need a way to make sure that higher priority workloads always get scheduled without being interrupted by lower priority workloads, this is where **Priority classes** come in.

* Priority classes help us define priorities for different workloads, so that higher priority workloads always get a priority over lower priority ones.
⇒ If a higher priority pod cannot be scheduled, the scheduler will try to eliminate a lower priority workload to make that happen. So **Priority classes** are non-namespace objects [meaning they're not created within a specific namespace, they're not attached to a specific namespace, so once they're created they're available to be configured on any pod in any namespace]

How to define priorities :-

↳ We'll define priorities with a range of numbers from -2,147,483,648 to 2,000,000,000. So a priority can be defined as high as one billion and as low as negative two billion (or) around that number.

1,000,000,000

↳ Applications send their results from the workloads in the form of requests to the system.

-2,147,483,648

higher priorities

↳ A larger number indicates higher priorities.

↳ The numbers that are mentioned above are for the applications that are deployed on the system.

↳ The numbers that are related to the workloads that are deployed on the system.

↳ But there is a separate range that is dedicated to the internal system components.

↳ These components have the highest priority of they're priorities.

↳ They're priorities are up to 2 billion.

2,000,000,000

↳ System

↳ These priorities are assigned to the system components.

1,000,000,000

↳ Apps

-2,147,483,648

List Priority Classes

To list priority class
in kubectl get priorityclass

How to create a new priority class

priority-class.yaml

apiVersion: scheduling.k8s.io/v1
kind: PriorityClass
metadata:
name: high-priority
value: 1000000000
description: "Priority class for mission critical pods"

[The "description" is an "optional field"] Once the "priority-class.yaml" definition file is created, we associate this priority class to a pod by using the "priorityClassName" property inside the "pod-definition".

Pod-definition.yaml

apiVersion: v1

kind: Pod

metadata:

name: nginx

labels:

name: nginx

spec:

containers:

- name: nginx

image: nginx

ports:

- containerPort: 8080

priorityClassName: high-priority

When that pod is created, it is assumed to have this priority that is assigned to the priority class

What if we don't set ~~a~~ priority class name for a pod?

⇒ By default, it is assumed to have a priority class value of "zero"

⇒ If you would like to modify that, you must create a "new priority class" & assign to global default property to "true"; This is an optional field but it can be used to define the default ~~priority~~ priority of all the pods

⇒ If there's no priority class name explicitly defined in that pod's configuration file

Global default property
NOTE: It can be defined only through "single priority class"

because we cannot have multiple default values

* we cannot have more than one priority class

with global default property set to true

Priority-class.yaml

apiVersion: scheduling.k8s.io/v1

kind: PriorityClass

metadata:

name: high-priority

value: 100000000

description: "Priority class for mission critical pods"

globalDefault: true

cer

Effect of Pod Priority

→ If there are two workloads with two priorities coming into be scheduled

* A critical app with a higher priority of 7
* And a jobs app with a lower priority of 5.

→ Since the "critical app" with the "higher priority" of "7" hence it gets placed "first" & if the no resource available, so the workload with the lower priority also gets placed next after the highest priority pod gets placed.

→ Now, we've a high priority job that comes in with a priority of 5 & there's no more resources available on the cluster, what happens?

→ This behavior is defined by the "preemption policy", defined in the "priority class" assigned to the new workload.

priority-class.yaml

apiVersion: scheduling.k8s.io/v1

kind: PriorityClass

metadata:

name: high-priority

(value: 100000000)

description: "Priority class for mission critical pods"

preemptionPolicy: PreemptLowerPriority # kills low priority pods

⇒ If the preemption policy is not set, its default value is set to "preempt lower priority". This means that it would kill existing lower priority job & take its place.

⇒ But if you do not want it to kill (or) evict the existing workload & instead wait for the cluster resources to free up then you must set this "preemption policy to never".

⇒ preemption policy to never makes the pods in that priority class non-preempting meaning that they cannot preempt other pods & it will wait to be scheduled in the scheduling queue. However, they will get a "higher priority" on scheduling over other lower priority pods that're also waiting to be scheduled.

⇒ So the priority still applies only when it is going to be scheduled, but will it preempt / kill / evict the existing workload or not this is defined by "the preemption policy".

You can compare the priority classes on both pods using the following command:-

Kubectl get pods -o custom-columns="NAME:.metadata.name,PRIORITY:.spec.priorityClassName"

Priority format: p0/p1/p2/p3