

81. Configuring profiles

How kubernetes Scheduler works

Consider scheduling a pod to one of the 4 nodes which're part of k8s cluster

pod-definition.yaml

apiVersion: v1

kind: Pod

metadata:

name: simple-webapp-color

spec:

priorityClassName: high-priority

containers:

- name: simple-webapp-color

image: simple-webapp-color

resources:

⇒ Using the above "pod-definition.yaml" a pod got created waiting to be scheduled in one of the 4 nodes in the cluster.

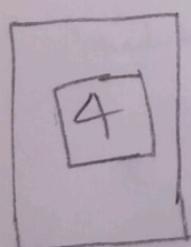
requests:

memory: "1Gi"

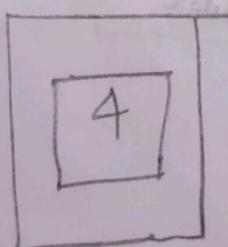
CPU: 10

⇒ It has a resource requirement of 10 CPU so it's gonna be scheduled on a node that has 10 CPU remaining

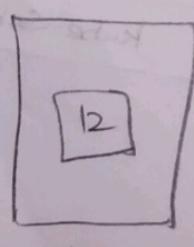
so see the available CPU on all of those nodes



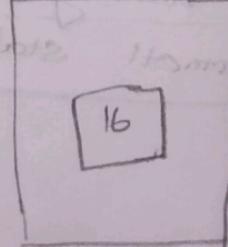
node 01



node 02



node 03



node 04

→ Now, the pod is not alone waiting to be scheduled. There're other pods too. Scheduling queue:
→ So, the first thing that happens is that, when these pods are created, the pods ends up in a "scheduling queue".
→ This is where the pods wait to be scheduled. So at this stage, pods are sorted based on the priority defined on the pods. In this case your pod has a high priority set.

To set a priority:

→ First, create a priority class that looks like this below. set it a "name" & set it a "priority value",
→ In this case, it is set to "1 million".

kind: PriorityClass

metadata:

name: high-priority

value: 1000000

globalDefault: false

description: "This priority class should be used for xyz service pods only"

"1 million" is really a very high priority, this is how the pod with very high priority gets to the beginning of the queue to be scheduled first, so that sorting happens in this scheduling phase.

• It's that another scenario is after that creation, to

Filtering: processes which do not have enough resources

→ Then, our pod enters the "filter" phase

→ This is where nodes that cannot run the pod are filtered out

→ In our case node01 & node02 do not have sufficient resources [do not have 10 CPU remaining] so they are filtered out

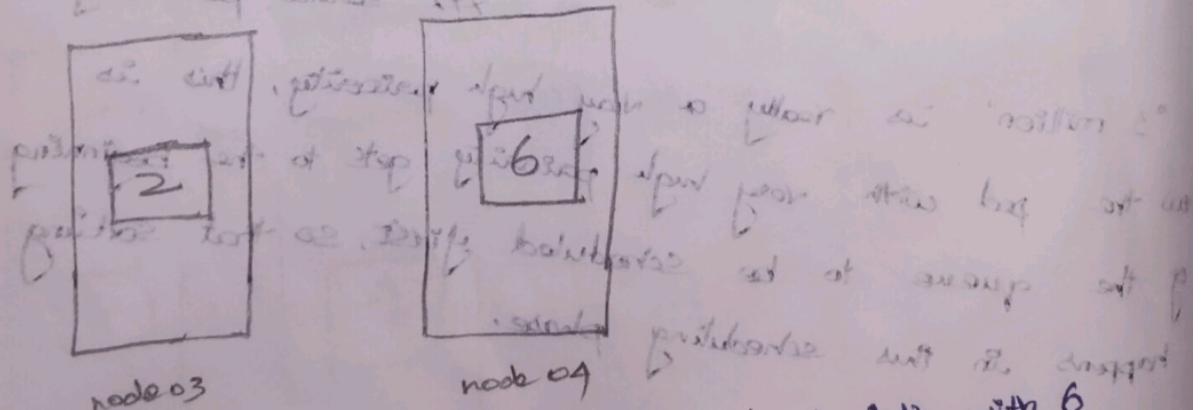
The next phase is of Scoring phase

Scoring phase:

→ In this phase, nodes are assigned with different weights

→ From the 2 remaining nodes node03 and node04, the "scheduler" associates a score to each node based on the free space that it will have after reserving the CPU required for that pod.

→ Currently node03 has CPU limit of 12, as the pod definition file it requires 10 CPU limit, so after reserving the CPU limit 10, it will have 2 similarly for node04 it will have 6 after reserving the CPU limit 10



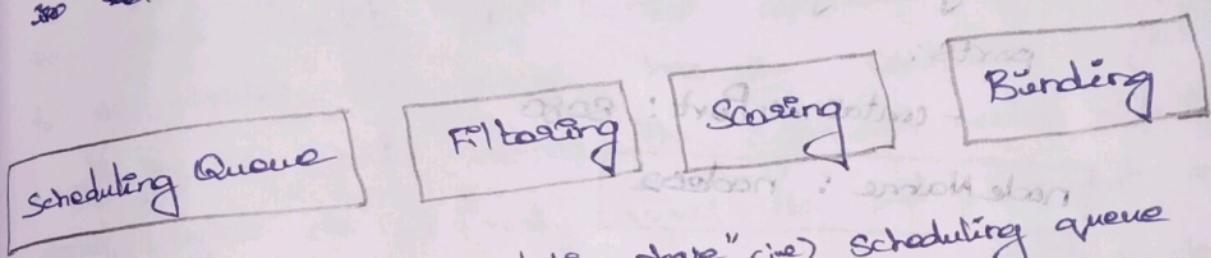
* node03 left with 2 whereas node04 left with 6

⇒ Since after reserving the pod to that node
 node 04 has maximum resources left, hence the
 "node 04" gets the highest score
 ⇒ Hence "node 04" gets picked up

Binding?

⇒ And finally in Binding phase, the pod is finally
 bound to a node with the highest score.

⇒ Now, all of the below operations are achieved
 via certain plugins



Eg: while in "scheduling phase" (i.e) scheduling queue
 its the "priority sort" plugin that sorts the pods in an
 order based on the priority configured on the pods,
 this is how pods with ~~higher priority~~ gets higher priority
 over the other pods when scheduling

In the Filtering Stage, its the NodeResourcesFilter
plugin that identifies the nodes that has "sufficient
 resources" required by the pods & filters out the nodes
 that doesn't

⇒ Other plugin in "filtering stage" - NodeName
 plugin that checks if a pod has a node name mentioned
 in the pods spec. if yes, it filters out all the nodes that
 does not match this name.

at last ["path distribution"]

drift and works between 0% and 100%
 collecting drift, "count" at last is "path distribution"

pod-definition.yaml provides info about

apiVersion: v1

kind: Pod

metadata:

name: nginx

labels:

name: nginx (select pod by following line)

spec:

containers:

- name: nginx

image: nginx

ports:

- containerPort: 8080

nodeName: node02

Node Unschedulable plugin:

kubectl describe node node02 controlplane

Name:

controlplane

Roles:

control-plane

CreationTime: Thu, 06 Oct 2020 06:19:51 -400

Taints:

node.kubernetes.io/unschedulable:NoSchedule

Unschedulable: "true"

Leave:

⇒ This plugin filters out nodes that has the

'unschedulable flag'

set to

"true"

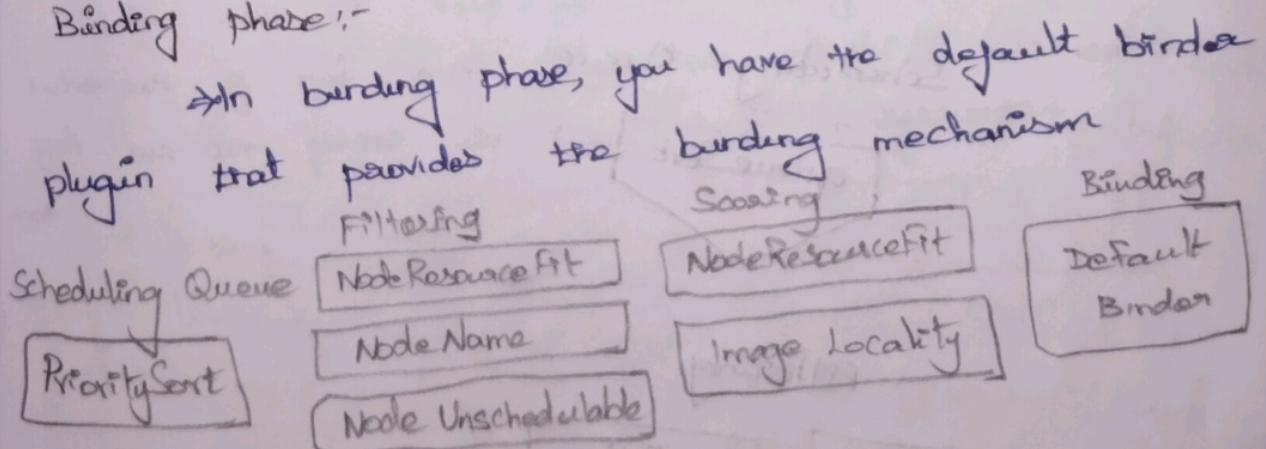
⇒ For the nodes, which has this

"unschedulable flag" is set to "true", this particular

plugin "Node Unschedulable" makes it impossible for pods to be scheduled on those nodes. In the scoring phase - NodeResourcesFit, the plugin associates a score to each node based on the resource available on it after the pod is allocated to it. So, as you can see, a single plugin can be associated in multiple different phases.

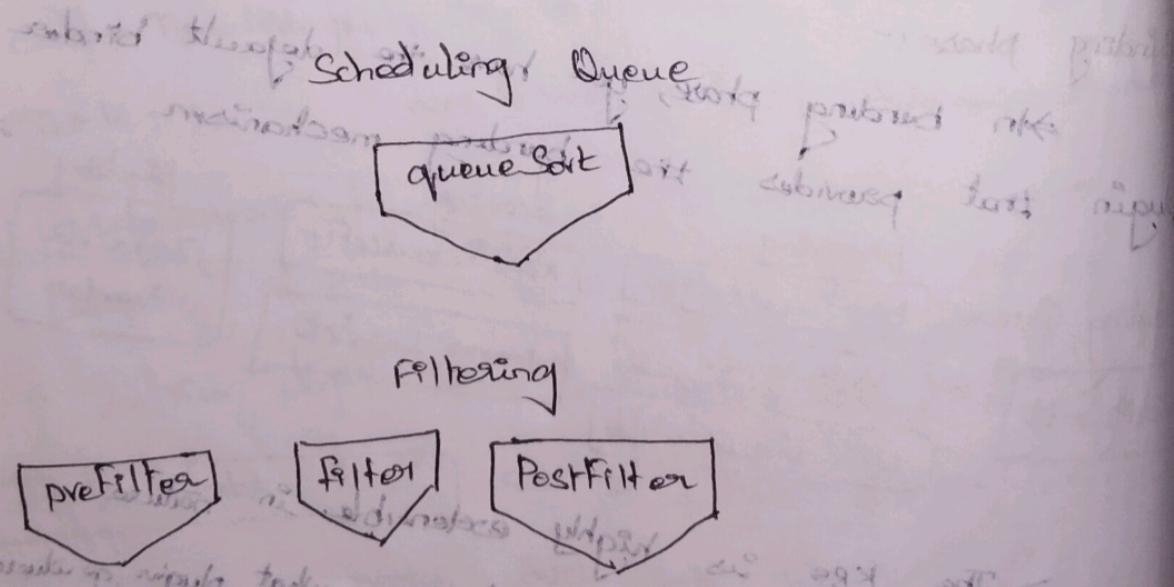
Another eg of plugin in this stage would be the image locality plugin that associates a high score to the nodes that already has the container image used by the pods among the different nodes at that place, plugins do not really reject the pod placement on a particular node, for e.g. - in case of image locality node, it ensures that pod are placed on a node that already has the image but if there are no nodes available, it'll anyway place the pod on a node that does not even have the image. so its just that Scoring that happens at this stage.

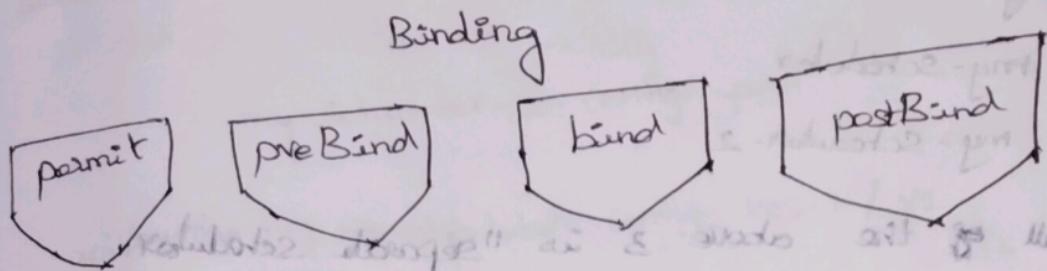
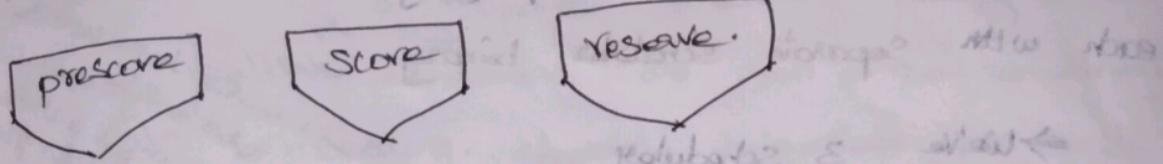
Binding phase:-



⇒ The K8s is highly extensible in nature makes it possible for us to customize what plugin goes where.

- for us to write our own plugin and plug them in here, that is achieved with the help of "extension points"
 - ⇒ At each stage, there is an extension point to which a plugin can be replugged to.
 - ⇒ In the scheduling queue, we've a queue sort extension to which the priority sort plugin is plugged to.
 - ⇒ And then we have the filter extension, that is score of the bind extension to which each of these plugins that we just talked about are plugged to.
 - ⇒ There're extensions before entering the filter phase, called "pre-filter extension" & after the filter phase, called post filter and then there are pre-score extensions before the score extension point or reverse after the extension point the score extension point point at last prebind task step etc.
 - ⇒ And there is permit of preBind before the bind & post bind after the "binding phase"





→ Basically we can get a custom code of your own to run anywhere in those points by just creating a plugin & plugging it into the respective kind of point that you want to plug it to. and making it

* Some additional plugins that come by default that are associated with the different extension points because some of the plugins span across multiple extension points

* Some of them are just within a specific extension point

Let's look at how we can change the default behavior of how these plugins are called & how we can get our own plugins in there if it's really needed.

122

Creating a custom reader plugin with the help of a subtask space in category algorithm. Dragged & dropped into workspace with the help of a subtask space & editor for test suite of certain problems from problems in current subtask space - things, strings etc.

→ Earlier, we talked about deploying 3 separate schedulers each with separate scheduler binary

→ We've 3 schedulers

1) Default Scheduler

2) my-scheduler

3) my-scheduler-2

1. config → API server
2. kubeconfig → config location
of binary & custom
Scheduler config file

Now all of the above 3 is "separate scheduler

binaries" that run with "separate scheduler config file" associated with each of them. This is one way of deploying "multiple schedulers".

The problem here is "since there are multiple separate processes there is an additional effort required to maintain those separate processes, & more importantly, since they are separate processes they may run into race conditions while making scheduling decision".

E.g.: One scheduler may schedule a workload on a node without knowing that there's another scheduler scheduling a workload on that same node at the same time.

→ So with the 1.18 release of k8s a feature to support multiple profiles in a single scheduler was introduced, so now you can configure multiple profiles within a single scheduler in the configuration file by adding more entries to the list of profiles & for each profile, specify a separate scheduler name.

So this creates a separate profile for each scheduler which acts as a "separate scheduler" itself except that now multiple schedulers are run in the same binary as opposed to creating separate binaries for each scheduler.

my-scheduler-2 config.yaml

apiVersion: kubescheduler.config.k8s.io/v1

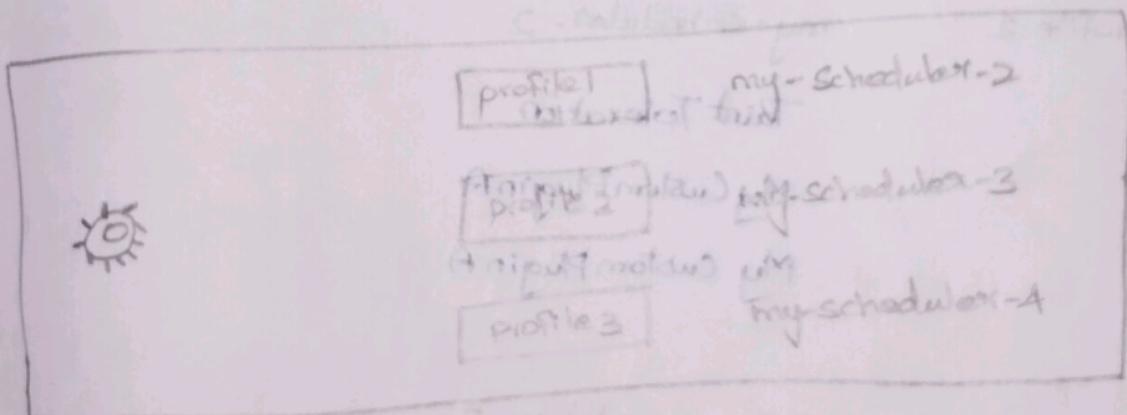
kind: KubeSchedulerConfiguration

profiles:

- schedulerName: my-scheduler-2

- schedulerName: my-scheduler-3

- schedulerName: my-scheduler-4



How to configure these different scheduler profile to work differently, because right now all of them has different names. So they are gonna work just like the "default Scheduler".

How do you configure them to work differently?

Under each "Scheduler profile" we can configure the plugins the way we want to.

For e.g.: for the "my-scheduler-2" profile, we're going to disable certain plugins like "taint tolerance" plugin & enable my own custom plugin.

→ For the "my-scheduler-3" profile, we're going to disable all the "pre score" and "score" plugins. So, this is how it is going to look, under the "plugins" section, specify the extension point and enable or disable the plugins by name or a pattern as shown -

profile 1 my-scheduler-2

Taint Tolerance

MyCustomPluginA

MyCustomPluginA

profile 2 my-scheduler-3

prescore score

profile 3 my-scheduler-4

my-scheduler-2-config.yaml

apiVersion: kubescheduler.config.k8s.io/v1

kind: KubeSchedulerConfiguration

profiles:

- schedulerName: my-scheduler-2

plugins:

score:

disabled:

- name: TaintToleration

enabled:

- name: MyCustomPluginA

- name: MyCustomPluginB

- schedulerName: my-scheduler-3

plugins:

preScore:

disabled:

- name: '*'

Score:

disabled:

- name: '*'

- schedulerName: my-scheduler-4

To read more about this checkout the k8s enhancement proposal that introduced multi scheduling profiles - [1451 - multi-scheduling-profiles](#)