

78. Multiple Schedulers I feel quite Tough 109
123 In this lecture, we look at deploying "multiple schedulers" in a k8s cluster

We already know how default scheduler works in k8s, it has an algorithm that distributes pods across nodes evenly as well as takes into considerations various conditions we specify through "labels & tolerations" & "node affinity" etc.,

But what if none of those satisfies your needs

say you've specific apps that require its components to be placed on nodes after performing some additional checks?

So you decide to have your own scheduling algorithm to place pods on nodes, so that you can add your own custom conditions & checks in it. k8s is highly extensible, we can write our own k8s scheduler program, package it & deploy it as the "default scheduler" or as an "additional scheduler" in the k8s cluster, that way all of the other apps can go through default scheduler.

However, some specific apps that you may choose can use your own custom scheduler

so, a k8s cluster can have "multiple schedulers" at a time

when creating a "pod" or "deployment", we can instruct k8s to have the pod scheduled by a "specific scheduler". Let see how to do this.

These are different schedulers & have different names so that we can identify them as separate schedulers, the "default scheduler" is named as "a default scheduler" & this name is configured in a "kube-scheduler configuration file" that look like this

Scheduler-config.yaml

```
apiVersion: kubescheduler.config.k8s.io/v1
```

```
kind: KubeSchedulerConfiguration
```

```
profiles:
```

```
- schedulerName: default-scheduler
```

Now, the default scheduler doesn't really need one because if you don't specify a name, it sets the name to a "default scheduler"

⇒ And for the other schedulers, we could create a separate configuration file. Set the "scheduler name" like this.

myScheduler-config.yaml

```
apiVersion: kubescheduler.config.k8s.io/v1
```

```
kind: KubeSchedulerConfiguration
```

```
profiles:
```

```
- schedulerName: my-scheduler
```

Deploy additional Scheduler Cluster runs on
Cluster runs on system to run schedules
→ Instead we saw how to "deploy the
"K8s Kube-Scheduler" as our scheduler under k8s
"we download the Kube-Scheduler binary & run
it as a service with a set of options
what <https://storage.googleapis.com/kubernetes-release/v1.12.0/bin/linux/amd64/kube-scheduler>

Now to deploy the additional scheduler, you may use
the same "Kube-Scheduler" binary (or) use one that
you might have built for yourself, which is what
you would do if you needed the scheduler to work
differently
Turn 2 pages before there
In this case, we'll use the same binary to deploy
the "additional scheduler". This time, we point the configuration
on to custom configuration file that we created, so
each scheduler uses a separate configuration file & with each
Kube-Scheduler service. file having its own scheduler name

ExecStart = /usr/local/bin/kube-scheduler
--config = /etc/kubernetes/config/kube-scheduler.yaml

my-scheduler-2.service
ExecStart = /usr/local/bin/kube-scheduler
--config = /etc/kubernetes/config/my-scheduler-2-config.yaml

And note there are other options to be passed in,
such as the kubeconfig file to authenticate into
K8s API but we're skipping this now to make it
Simple.

This is not how we deploy a custom scheduler because without "kubeadm" deployment all the control plane components run as a daemonset and deployment within the k8s cluster.

We want to follow this Systemd binary service approach method 99% of time.

Deploy Scheduler as a pod

Deploy additional scheduler as pod

we have to create a "pod-definition file" & specify the "kube config" property which is the path to the "schedulerconfig" file that has the authentication information to connect to k8s API server.

we have to pass in our custom configuration file as a config option to the scheduler configuration file as path to scheduler config file [it has the authentication information to connect to k8s API server].

Note: config = path to our custom scheduler configuration file.

we've the scheduler name specified in the file, so that's how the name gets picked up by the scheduler.

my-custom-scheduler.yaml

apiVersion: v1
kind: Pod

metadata:

name: my-custom-scheduler
namespace: kube-system

spec:
containers:

- command:

- kube-scheduler

- address = 127.0.0.1

- --kubeconfig = /etc/kubernetes/scheduler.conf

- --config = /etc/kubernetes/my-scheduler-config.yaml

image: k8s.gcr.io/kube-scheduler-amd64:v1.11.3

name: kube-scheduler

→ my-scheduler - config.yaml

apiVersion: kubescheduler.config.k8s.io/v1

kind: kubeschedulerConfiguration

profiles:

- schedulerName: my-scheduler

leaderElection:

leaderElect: true

resourceNamespace: kube-system

resourceName: lock-object-my-scheduler

⇒ Now, another important option to look here is the "leaderElect" option, if this goes into the "kubescheduler" configuration.

⇒ The "leaderElect" option is used, when you have multiple copies of the scheduler running on different master nodes or on high-availability setup, where you have multiple master nodes with the k8s scheduler process running on both of them.

⇒ If multiple copies of the same scheduler are running on different nodes, only one can be active at a time & that's where the "leaderElect" option helps in choosing a leader who will lead the scheduling activities.

Remember:-

If you have multiple masters, you can pass in this additional parameter to set a log object name & this is to differentiate the new custom scheduler from the default election process.

Deploying additional scheduler as a deployment

Now we can see the new custom scheduling running.

Kubelet gets pods -- namespace = kube-system

kube-scheduler-master	1/1	Running	0	1m
my-custom-scheduler	1/1	Running	0	9s

Use Custom Scheduler

Once, we have deployed that custom scheduler, the next step is to configure a pod or a deployment to use this "new scheduler".

How do we use our custom scheduler? In the pod-definition file, add a new field called "schedulerName" & specify the name of the "new scheduler". This way when the pod is created, the right scheduler gets picked up & the scheduling process works.

Pod-definition.yaml

kind: Pod

metadata:

name: nginx

spec:

- containers:

- image: nginx

name: nginx

schedulerName: my-custom-scheduler

Now create the pod

`kubectl create -f pod-definition.yaml`

→ If the Scheduler was not configured correctly, then the pod will continue to remain in Pending State, if everything is good, then the pod will be in Running State.

→ If the pod is in pending state, then we can take the log under the pod describe command.

→ kubectl describe command & you'll mostly notice that the scheduler isn't configured correctly.

View Events :-
How do you know which scheduler picked up a particular pod
If we've multiple schedulers, we can view this in the events using the

`kubectl get events -o wide`

→ This will list all the events in the current namespace look for "scheduled" events & you can see the source of the events is the custom scheduler that we created.
Y100X Events → that's the name that we give to the custom scheduler & the message - successfully assigned default/
source to pod

LASTSEEN	COUNT	NAME	KIND	TYPE	REASON	SOURCE	MESSAGE
9s	1	nginx-15	Pod	Normal	Scheduled	my-custom-scheduler	successfully assigned default/nginx to workerd

View Scheduler Logs

In case if you run into issues, for that, view the logs using the `kubectl logs` command & provide the scheduler name either the "Pod name" or "deployment name" & then "right namespace".

`kubectl tail my-custom-scheduler -n kube-system`

my-custom-scheduler is the log output from the pod.

git clone https://github.com/kube-scheduler/test-with-kubelet.git

cd kubelet/test-with-kubelet

./start.sh

curl -H "Content-Type: application/json" -X POST -d '{ "apiVersion": "kubescheduler.config.k8s.io/v1", "kind": "KubeSchedulerConfiguration", "clientConnection": { "kubeConfig": "/etc/kubernetes/config/scheduler.conf" }, "leaderElection": { "leaderElect": false } }' http://127.0.0.1:10250/api/v1/nodes/127.0.0.1/schedule

curl -H "Content-Type: application/json" -X POST -d '{ "apiVersion": "kubescheduler.config.k8s.io/v1", "kind": "KubeSchedulerConfiguration", "clientConnection": { "kubeConfig": "/etc/kubernetes/config/scheduler.conf" }, "leaderElection": { "leaderElect": false } }' http://127.0.0.1:10250/api/v1/nodes/127.0.0.1/schedule

curl -H "Content-Type: application/json" -X POST -d '{ "apiVersion": "kubescheduler.config.k8s.io/v1", "kind": "KubeSchedulerConfiguration", "clientConnection": { "kubeConfig": "/etc/kubernetes/config/scheduler.conf" }, "leaderElection": { "leaderElect": false } }' http://127.0.0.1:10250/api/v1/nodes/127.0.0.1/schedule

curl -H "Content-Type: application/json" -X POST -d '{ "apiVersion": "kubescheduler.config.k8s.io/v1", "kind": "KubeSchedulerConfiguration", "clientConnection": { "kubeConfig": "/etc/kubernetes/config/scheduler.conf" }, "leaderElection": { "leaderElect": false } }' http://127.0.0.1:10250/api/v1/nodes/127.0.0.1/schedule

How to access & modify the Kube-Scheduler-Service, if your cluster uses Systemd to run the Scheduler instead of static pods [which is rare in hub-and-spoke clusters but can occur in custom setup or learning environment].

Step by step - Modifying Kube-Scheduler-Service

This applies only when your cluster is running Kube-Scheduler as Systemd Service, not as a static pod in /etc/kubernetes/manifest.

Step 1: Locate the Kube-Scheduler Service

check if its managed by Systemd

Systemctl status kube-scheduler

You will see output like

• Kube-Scheduler Service - Kubernetes scheduler
loaded: loaded (/etc/systemd/system/kube-scheduler.service)
enabled: enabled (disabled by kubelet)

Active: active (running)

Step 2: Edit the Systemd unit

Sudo systemctl edit kube-scheduler

This opens an override file you'll usually see:

/etc/systemd/system/kube-scheduler.service

[Service]

ExecStart=

/usr/bin/kube-scheduler --config=/etc/kubernetes/

-key/config/kubescheduler.yaml

The first ExecStart= line closes the default so the new one overrides it.

[Make sure the config file path of binary match your environment]

Step 3: Reload and Restart

Sudo systemctl daemon-reexec

Sudo systemctl daemon-reload kubelet

Sudo systemctl restart kube-scheduler

Confirm its running:-

Systemctl status kube-Scheduler