

# KUBERNETES

## Scheduling & Resources



# Manual Scheduling

- By default, nodeName is not set. Kubernetes adds it automatically.
- The scheduler goes through all the PODs and looks for those that don't have this set.
  - Identify the Node.
  - Schedule the POD on the Node.
  - Create 'Bind Object': bind the POD to the Node.
- The PODs will stay in 'Pending' status if there is no scheduler.
- Create a binding object and send a post request to the POD binging API that mimics the actual scheduler.

```
1  apiVersion: v1
2  kind: Binding
3  metadata:
4    name: nginx
5  target:
6    apiVersion: v1
7    kind: Node
8    name: <Node-Name>
```

# Labels and Selectors

- Labels: properties attached to each item.
- Selectors: helps to filter the labels.

## How?

- In the POD definition file

- You can search from label:

```
$ kubectl <get-pods> --selector <label-type>=<label>
```

```
1  apiVersion: v1 # String
2  kind: Pod # String
3  metadata: # Dictionary
4    name: my-app
5  labels: # Dictionary
6    app: my-app
7    type: front-end
```

# Taints and Tolerations

- Taints and Tolerations are used to set restrictions on what PODs can be scheduled on a Node.
- To taint a Node:  
**\$ kubectl taint nodes <node-name> <key>=<value>: <taint-effect>**

## 3 taint effects:

- NoSchedule
- PreferNoSchedule
- NoExecute
- See all taints Nodes:  
**\$ kubectl describe node kubemaster | grep Taint**

# NodeSelector

- If we want a small a small POD will go to a small Node and a big POD will go to a big Node.
- We label them in the POD definition file with NodeSelector

```
apiVersion: v1 # String
kind: Pod # String
metadata: # Dictionary
  name: my-app
  labels: # Dictionary
    app: my-app
    type: front-end
spec:
  containers: # List of Dictionaries
    - name: my-app-container # String
      image: myregistry.com/my-app:1.0 # String

  # ----- NODE SELECTOR ----- #
  nodeSelector: # List of
    size: Large # String
```

# NodeAffinity

- Provide us with advanced capabilities to limit POD placement on specific Nodes.
- Makes sure PODs are hosted on a Node.

```
1  apiVersion: v1 # String
2  kind: Pod # String
3  metadata: # Dictionary
4    name: my-app
5    labels: # Dictionary
6      app: my-app
7      type: front-end
8  spec:
9    containers: # List of Dictionaries
10   - name: my-app-container # String
11     image: myregistry.com/my-app:1.0 # String
12   # ----- NODE AFFINITY -----
13   affinity:
14     requiredDuringSchedulingIgnoredDuringExecution:
15       nodeSelector: # Dictionary
16         - matchExpressions:
17           - key: size
18             operator: In
19             values:
20               - Large
```

# 2 Types of NodeAffinity.

## 1. Available:

- requiredDuringSchedulingIgnoredDuringExecution
- preferredDuringSchedulingIgnoredDuringExecution

## 2. Planned:

- requiredDuringSchedulingRequiredDuringExecution

## NodeAffinity vs Taints & Tolerations

- Each alone can't guarantee clean POD to Node.
- Both together can make sure clean POD to Node



# Resource Requirement and Limits

- ‘kube-scheduler’ decides what POD goes to each Node.
- If there is no available resources, the POD will be in ‘Pending’ state.
- You can specify the amount of resources in the POD-Definition.yaml file

```
apiVersion: v1 # String
kind: Pod # String
metadata: # Dictionary
  name: my-app
  labels: # Dictionary
    app: my-app
    type: front-end
spec:
  containers: # List of Dictionaries
    - name: my-app-container # String
      image: myregistry.com/my-app:1.0 # String
      # ----- RESOURCES -----
      resources:
        requests:
          memory: "1Gi"
          cpu: 1
        limits:
          memory: "2Gi"
          cpu: 2
```

# Limit Ranges

- Can help to define default values to be set for containers in PODs that are created without a request or limit specified in the ‘POD-Definition.yaml’

```
1 apiVersion: v1
2 kind: LimitRange
3 metadata:
4   name: cpu-resource-constraint
5 spec:
6   limits:
7     - default:
8       cpu: 500m # Limit
9       defaultRequest:
10      cpu: 500m # Request
11      max:
12      cpu: "1" # Limit
13      min:
14      cpu: 100m # Request
15      type: Container
16
```

# DaemonSets

- Like ReplicaSets, DaemonSets helps you to deploy multiple instances of PODs.
  - Runs each copy of a POD on each Node
- Whenever a new Node is added to the cluster, a replica of the POD is auto-added to the Node. When a Node is removed, the POD is removed as well.
- DaemonSets makes sure the POD is in each Node in the Cluster.

## Uses:

- Monitor agent
- Log collector
- Networking

## DaemonSets-definition.yaml

```
1 apiVersion: apps/v1
2 kind: DaemonSet
3 metadata:
4   name: monitoring-Daemon
5 spec:
6   selector:
7     matchLabels:
8       app: monitoring-agent
9   template:
10    metadata:
11      labels:
12        app: monitoring-agent
13    spec:
14      containers:
15        - name: monitoring-agent
16          image: <your-docker-image-registry>/<your-docker-image-name>:latest
17
```

# Static PODs

- If there is no: master, cluster, or even Nodes can the POD stand alone?
- Without kube-apiserver, the kublet is configured to read the POD.yaml files from a directory in a server that stores all the information about the PODs.
- The kublet checks the directory for changes: edit/add/delete.

Change the directory the kublet checks:

- In kublet.service:  
**--pod-manifest-path=/path/to/dir \\**
- In kubeconfig.yaml:  
**staticPodPath:/path/to/dir**
- The API server can create both static PODs and Cluster Nodes and can notice the changes in the static Nodes.

# Multiple Schedulers

- We can make a custom scheduler and also make it the default scheduler.

```
1 apiVersion: kubernetes.config.k8s.io/v1
2 kind: KubeSchedulerConfiguration
3 profiles:
4   - schedulerName: Custom-Scheduler
```

- Also, we can add an additional scheduler by editing ‘my-scheduler.service’

**--config-/etc/kubernetes/config/<my-scheduler.yaml>**

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: my-custom-scheduler
5   namespace: kube-system
6 spec:
7   containers:
8     command:
9       - kube-scheduler
10      - --address=127.0.0.1
11      - --kubeconfig=/path/to/my-custom-scheduler
12      - --config=/path/to/my-custom-scheduler.yaml
13
14     image: k8s.gcr.io/kube-scheduler-amd64:v1.11.3
15     name: kube-scheduler
```

# Schedule Profiles

- Schedule Queue: The place where the PODs are waiting to be created.
  - PODs are stored based on priority.
- Filter: if there are any resources to use
- Scoring: Nodes are scored in different weights– based on a free space.
- Binding: the processes of a POD bind to a Node.
- Also, we can add an additional scheduler by editing ‘my-scheduler.service’

# Scheduling Plugins

## Scheduling Queue:

- Queuesort
- PrioritySort

## Filtering:

- NodeResourcesFit
- NodeName
- NodeUnschedulable
- NodeResourcesFit
- TaintToleration
- NodePorts
- NodeAffinity

## Scoring:

- NodeResourceFit
- ImageLocality

## Binding:

DefaultBinder