Dynamically rebalancing workloads and optimizing resource utilization in Kubernetes

Chris Nesbitt-Smith





Chris Nesbitt-Smith

UK Gov | LearnK8s | Control Plane | lots of open source



Datacentre as a single computer























Combining autoscalers













Downscaling and fragmentation











```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-deployment
spec:
  replicas: 5
  selector:
    matchLabels:
      app: nginx
  template: ←
                                         pod definition
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx:1.14.2
```



```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-deployment
spec:
  replicas: 5
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx:1.14.2
```

There's no field for "rebalance"!

Descheduler





OBO Defrag V8 Professional Edition					
Elle View Defragmentation Jobs Reports 2					
Q ♀ ⊯・ II ⊠ ₽ ₽ ₽ ₽ ₽ ₽ ₽	9 0	0			
Volume Current File/Folder	Status	Total Files	Fr S	ize [M8]	
	100%	98896	20	127992	
	0.76	0		00470	
Cluster C Jobs Reports					
C: - 4750 Clusters per Block					100 %
Ready					









Kubernetes scheduler









Scheduling

Binding



PodFitsHostPorts PodFitsHost **PodFitsResources PodMatchNodeSelector NoVolumeZoneConflict** NoDiskConflict MaxCSIVolumeCount CheckNodeMemPressure CheckNodePIDPressure CheckNodeDiskPressure CheckNodeCondition

PodToleratesNodeTaints

CheckVolumeBinding



SelectorSpreadPriority

InterPodAffinityPriority

LeastRequestedPriority

MostRequestedPriority

RequestedToCapacityRatioPriority

BalancedResourceAllocation

NodePreferAvoidPodsPriority

NodeAffinityPriority

TaintTolerationPriority

ImageLocalityPriority

ServiceSpreadingPriority

EqualPriority

EvenPodsSpreadPriority



Descheduler policies

Descheduler plugins

Name	Extension Point Implemented	Description
RemoveDuplicates	Balance	Spreads replicas
LowNodeUtilization	Balance	Spreads pods according to pods resource requests and node resources available
HighNodeUtilization	Balance	Spreads pods according to pods resource requests and node resources available
RemovePodsViolatingInterPodAntiAffinity	Deschedule	Evicts pods violating pod anti affinity
RemovePodsViolatingNodeAffinity	Deschedule	Evicts pods violating node affinity
RemovePodsViolatingNodeTaints	Deschedule	Evicts pods violating node taints
RemovePodsViolatingTopologySpreadConstraint	Balance	Evicts pods violating TopologySpreadConstraints
RemovePodsHavingTooManyRestarts	Deschedule	Evicts pods having too many restarts
PodLifeTime	Deschedule	Evicts pods that have exceeded a specified age limit
RemoveFailedPods	Deschedule	Evicts pods with certain failed reasons

$\bullet \bullet \bullet$

apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 pluginConfig:
 - name: "RemoveDuplicates"
 - name: "RemovePodsHavingTooManyRestarts"
 args:

Almost a CRDI

podRestartThreshold: 100

includingInitContainers: true

plugins:

deschedule:

enabled:

- "RemovePodsHavingTooManyRestarts"

balance:

enabled:

- "RemoveDuplicates"






apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 pluginConfig:
 - name: "RemoveDuplicates"
 - name: "RemovePodsHavingTooManyRestarts"
 args:
 - podRestartThreshold: 100
 - includingInitContainers: true
 - plugins.
 - deschedule:
 - enabled:
 - "RemovePodsHavingTooManyRestarts"
 - balance:
 - enabled:
 - "RemoveDuplicates"

pluging

apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 pluginConfig:
 - name: "RemoveDuplicates"
 - name: "RemovePodsHavingTooManyRestarts"
 args:

podRestartThreshold: 100

includingInitContainers: true

plugins:

deschedule: ←

enabled:

- "RemovePodsHavingTooManyRestarts"

extension points

balance: ← enabled:

- "RemoveDuplicates"





Extension points



Restart policy

```
apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:
```

- name: ProfileName
 pluginConfig:
 - name: "PodLifeTime"
 args:

```
maxPodLifeTimeSeconds: 10
```

- plugins:
 - deschedule:
 - enabled:
 - "PodLifeTime"

pods older than 10 seconds are removed





Demo



Descheduler deployment



Descheduler deployment

Job CronJob Deployment



Descheduler deployment

Job CronJob Deployment



Descheduler deployment

Job
 CronJob
 Deployment



CronJob



```
apiVersion: batch/v1
kind: CronJob
metadata:
  name: descheduler-cronjob
  namespace: kube-system
spec:
  schedule: "*/1 * * * *"
  concurrencyPolicy: "Forbid"
  jobTemplate:
    spec:
      template:
        metadata:
          name: descheduler-pod
        spec:
          containers:
            - name: descheduler
              image: registry.k8s.io/desch...
```

Frequency













Deployment

















Duplicate policy







Ļ









Demo



Collecting metrics













Kubernetes API server



High utilization policy

apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 pluginConfig:
 - name: "HighNodeUtilization"

args:

thresholds:

"memory": 20

plugins:

balance:

enabled:

threshold

- "HighNodeUtilization"







underutilized appropriately utilized node




Allocatable























HOW TO (RIGHT) SIZE YOUR KUBERNETES CLUSTER FOR EFFICIENCY



Low utilization policy

apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 - pluginConfig:
 - name: "LowNodeUtilization"
 - args:













Demo



Node Problem Detector





























Node controller











node.kubernetes.io/unreachable





Node Problem Detector + Descheduler + Cluster Autoscaler







The Node Problem detector pod detects that the node is unreachable







Taints violation policy



apiVersion: "descheduler/v1alpha2"
kind: "DeschedulerPolicy"
profiles:

- name: ProfileName
 pluginConfig:
 - name: "RemovePodsViolatingNodeTaints"
 plugins:
 - deschedule:
 - enabled:
 - "RemovePodsViolatingNodeTaints"



The descheduler evicts all pods from the node

node.kubernetes.io/unreachable









Cluster Autoscaler





Takeaways

Recap


2. Descheduler

3. Descheduler policies

4. Metrics pipeline



2. Descheduler

3. Descheduler policies

4. Metrics pipeline



2. Descheduler

3. Descheduler policies

4. Metrics pipeline



2. Descheduler

3. Descheduler policies

4. Metrics pipeline



2. Descheduler

3. Descheduler policies

4. Metrics pipeline



Thank you!

Chris Nesbitt-Smith cns.me





Chris Nesbitt-Smith cns.me



Requests & limits



0.25 vCPU

512MB



256MB

Requests → scheduler Limits → kubelet QoS

Scheduler hints



Karpenter



Guaranteed, burstable, best-effort

















