

Chapter 2.2.2 – Deployment, Jobs and scaling

The screenshot shows a Windows desktop environment with a browser window open to a Qwiklabs course session. The title bar reads 'Deployments 1 | Qwiklabs'. The main content area is titled 'Deployment usage' and contains four icons with corresponding text descriptions:

- Roll out updates to the Pods** (Icon: Two overlapping arrows pointing right)
- Roll back Pods to previous revision** (Icon: Two overlapping arrows pointing left)
- Scale or autoscale Pods** (Icon: Three arrows forming a triangle)
- Well-suited for stateless applications** (Icon: A lightning bolt)

A callout box at the bottom center says 'applications stateless'. On the left sidebar, under 'Video Deployments 1', there is a list of topics including 'Ways to Create Deployments', 'Services and Scaling', 'Quiz: Deployments', 'Updating Deployments', 'Rolling Updates', 'Blue-Green Deployments', 'Canary Deployments', and 'Quiz: Update'. A 'Chat' button is located in the bottom right corner.

The screenshot shows a Windows desktop environment with a browser window open to a Qwiklabs course session. The title bar reads 'Deployments 1 | Qwiklabs'. The main content area is titled 'Deployment is a two-part process' and features a flow diagram:

```
graph TD; A[.yaml file] --> B[Deployment object]; B --> C[Deployment controller]; C --> D[Node]; D --> E["Pod Pod Pod"];
```

A callout box at the bottom says 'during this process a replica set is created'. To the right of the diagram, a man in a grey blazer is gesturing while speaking. On the left sidebar, under 'Video Deployments 1', there is a list of topics including 'Ways to Create Deployments', 'Services and Scaling', 'Quiz: Deployments', 'Updating Deployments', 'Rolling Updates', 'Blue-Green Deployments', 'Canary Deployments', and 'Quiz: Update'. A 'Chat' button is located in the bottom right corner.

Deployments 1 | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49529

Mark as Completed

Depployments 1

Deployment object file in YAML format

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: my-app
spec:
  replicas: 3
  template:
    metadata:
      labels:
        app: my-app
    spec:
      containers:
        - name: my-app
          image: gcr.io/demo/my-app:1.0
          ports:
            - containerPort: 8080
```

object file in emal format
the deployment named

Chat

Video Deployments 1

Video Ways to Create Deployments

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating

Deployments 1 | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49529

Mark as Completed

Depployments 1

Deployment has three different lifecycle states

are running
finally the failed state occurs when

Chat

Video Deployments 1

Video Ways to Create Deployments

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating

Ways to Create Deployments

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49530

Ways to Create Deployments

Mark as Completed

1 `$ kubectl apply -f [DEPLOYMENT_FILE]`

2 `$ kubectl run [DEPLOYMENT_NAME] \ --image [IMAGE]:[TAG] \ --replicas 3 \ --Labels [KEY]=[VALUE] \ --port 8080 \ --generator deployment/apps.v1 \ --save-config`

deployment imperative Li using a cube
CTL run command

Chat

Ways to Create Deployments

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49530

Ways to Create Deployments

Mark as Completed

3

Create a deployment

A deployment is a configuration which defines how Kubernetes deploys, manages, and scales your container image. Kubernetes will ensure your system matches this configuration.

Deployment

Container

Container image: nginx:latest
Select Google Container Registry image

Environment variables

Initial command (optional)

Done Cancel

Application name: nginx-1
Namespace: default
Labels: app: nginx-1
+ Add label

Cluster

Your deployment will use compute instances managed in a logical grouping called a 'cluster', which will be configured in a way that's great for getting started with Kubernetes.
The cluster will be named nginx-1-cluster
Zone: us-central1-a

Deploy View YAML

option is to use the gke were close menu in the GCP console here you can

Chat

Ways to Create Deployments

Ways to Create Deployments

Or output the Deployment config in a YAML format

```
$ kubectl get deployment [DEPLOYMENT_NAME]
```

```
master $ kubectl get deployment nginx-deployment
NAME      DESIRED   CURRENT   UP-TO-DATE   AVAILABLE   AGE
nginx-deployment   3          3          3           3          3m
```

```
$ kubectl get deployment [DEPLOYMENT_NAME] -o yaml > this.yaml
```

also output the deployment configuration
in a Yaml format this

Chat

Ways to Create Deployments

Ways to Create Deployments

Use the 'describe' command to get detailed info

```
$ kubectl describe deployment [DEPLOYMENT_NAME]
```

```
master $ kubectl describe deployment nginx-deployment
Name:           nginx-deployment
Namespace:      default
CreationTimestamp: Fri, 12 Oct 2018 15:23:46 +0000
Labels:          app=nginx
Annotations:    deployment.kubernetes.io/revision=1
Selector:        app=nginx
Replicas:       3 desired | 3 updated | 3 total | 3 available | 0 unavailable
StrategyType:   RollingUpdate
MinReadySeconds: 0
RollingUpdateStrategy: 25% max unavailable, 25% max surge
Pod Template:
  Labels:  app=nginx
  Containers:
    nginx:
      Image:  nginx:1.15.4
      Port:   80/TCP
      Host Port: 80/TCP
```

deployment use the cube CTL describe
command you'll learn

Chat

Ways to Create Deployments

Or use the GCP Console

nginx-1

Revision	Name	Status	Summary	Created on	Pods running/Pods total
1	nginx-1-7cb5b65464-2j97f	OK	nginx: nginx latest	Oct 12, 2018, 11:08:29 AM	3/3

Managed pods

Revision	Name	Status	Restarts	Created on
1	nginx-1-7cb5b65464-2j97f	Running	0	Oct 12, 2018, 11:08:29 AM
1	nginx-1-7cb5b65464-gct8h	Running	0	Oct 12, 2018, 11:08:29 AM
1	nginx-1-7cb5b65464-tz85f	Running	0	Oct 12, 2018, 11:08:29 AM

here you can see detailed information about the deployment revision

Chat

Services and Scaling

Scaling a Deployment manually

\$ kubectl scale deployment [DEPLOYMENT_NAME] --replicas=5

ACTIONS

Autoscale
Expose
Rolling Update
Scale

Scale

Replicas: 1

however at some point you'll probably need to scale the deployment

Chat

Services and Scaling | Qwikl 5:25 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49531

Services and Scaling

Mark as Completed

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Autoscaling a Deployment

```
$ kubectl autoscale deployment [DEPLOYMENT_NAME] --min=5 --max=15 --cpu-percent=75
```

Autoscale
Automatically scale the number of pods.

Minimum number of Pods (Optional):

Maximum number of Pods:

Target CPU utilization in percent (Optional):

Subtitles/closed captions (c)

Chat

by specifying the minimum and maximum



Services and Scaling | Qwikl 5:26 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49531

Services and Scaling

Mark as Completed

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Thrashing is a problem with any type of autoscaling

Cooldown/delay support:

```
--horizontal-pod-autoscaler-downscale-stabilization
```

deployed replicas frequently fluctuates because the metric

Subtitles/closed captions (c)

Chat

Quiz: Deployments | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49532

Quiz: Deployments

Your score: 100% Passing score: 50%

Congratulations! You passed this assessment.

Retake

Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google

✓ 1. What type of application is suited for use with a Deployment?

Batch

Written in Go

Stateful

Stateless

That is correct.

✓ 2. What is the relationship between Deployments and ReplicaSets?

There is no relationship; in modern Kubernetes, Replication Controllers are typically used to maintain a set of Pods in a running state.

A ReplicaSet configures a Deployment controller to create and maintain a specific version of the Pods that the Deployment specifies.

A Deployment configures a ReplicaSet controller to create and maintain all the Pods that the Deployment specifies, regardless of their version.

A Deployment configures a ReplicaSet controller to create and maintain a specific version of the Pods that the Deployment specifies.

There is no relationship; in modern Kubernetes, Replication Controllers are typically used to maintain a set of Pods in a running state.

Chat

Quiz: Deployments | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49532

Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google

✓ 2. What is the relationship between Deployments and ReplicaSets?

There is no relationship; in modern Kubernetes, Replication Controllers are typically used to maintain a set of Pods in a running state.

A ReplicaSet configures a Deployment controller to create and maintain a specific version of the Pods that the Deployment specifies.

A Deployment configures a ReplicaSet controller to create and maintain all the Pods that the Deployment specifies, regardless of their version.

A Deployment configures a ReplicaSet controller to create and maintain a specific version of the Pods that the Deployment specifies.

That is correct.

Chat

Updating Deployments | Q

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49533

Updating Deployments

Mark as Completed

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Updating a Deployment

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: my-app
spec:
  replicas: 3
  template:
    spec:
      containers:
        - name: my-app
          image: gcr.io/demo/my-app:1.0
          ports:
            - containerPort: 8080
```

```
$ kubectl apply -f [DEPLOYMENT_FILE]
```

```
$ kubectl set image deployment
[DEPLOYMENT_NAME] [IMAGE] [IMAGE]:[TAG]
```

command this allows you to change the pod template



Chat

Updating Deployments | Q

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49533

Updating Deployments

Mark as Completed

Video Services and Scaling

Quiz Quiz: Deployments

Video Updating Deployments

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Updating a Deployment

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: my-app
spec:
  replicas: 3
  template:
    spec:
      containers:
        - name: my-app
          image: gcr.io/demo/my-app:1.0
          ports:
            - containerPort: 8080
```

```
$ kubectl edit \
deployment/[DEPLOYMENT_NAME]
```

this opens the specification file using the vim editor that allows you to make



Chat

Updating Deployments | Q

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49533

Updating Deployments

Mark as Completed

- ✓ Video Services and Scaling
- ✓ Quiz Quiz: Deployments
- ✓ Video Updating Deployments**
- ✓ Video Rolling Updates
- ✓ Video Blue-Green Deployments
- ✓ Video Canary Deployments
- ✓ Quiz Quiz: Updating deployments
- ✓ Video Managing Deployments
- ✓ Video Lab Intro

Updating a Deployment

REFRESH EDIT DELETE ACTIONS KUBECTL

Rolling update

Update workload Pods to a new application version.

Minimum seconds ready (Optional): 0

Maximum surge (Optional): 30%

Maximum unavailable (Optional): 1

Container name Image

nginx nginx:latest

from the GCB console and perform a rolling update along

Chat

Updating Deployments | Q

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49533

Updating Deployments

Mark as Completed

- ✓ Video Services and Scaling
- ✓ Quiz Quiz: Deployments
- ✓ Video Updating Deployments**
- ✓ Video Rolling Updates
- ✓ Video Blue-Green Deployments
- ✓ Video Canary Deployments
- ✓ Quiz Quiz: Updating deployments
- ✓ Video Managing Deployments
- ✓ Video Lab Intro

The process behind updating a Deployment

Deployment

OLD ReplicaSet
Pod

NEW ReplicaSet
Pod Pod Pod

set this is an example of a rolling update strategy also known as a

Chat

Rolling Updates | Qwiklabs 5:35 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Set parameters for your rolling update strategy

Deployment

Max unavailable = 25%

Old ReplicaSet

means you want to have at least 75 percent

Chat

Rolling Updates | Qwiklabs 5:35 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Set parameters for your rolling update strategy

Deployment

Max surge = 2

New ReplicaSet

set for example you can specify

Settings

Chat

Rolling Updates | Qwiklabs 5:35 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Set parameters for your rolling update strategy

Deployment

Max surge = 25%

Old ReplicaSet

Pod Pod

New ReplicaSet

Pod Pod
Pod

the deployment controller looks at the total number

Settings Chat

Rolling Updates | Qwiklabs 5:36 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

An example of a rolling update strategy

```
[...]
kind: deployment
spec:
  replicas: 10
  strategy:
    type: RollingUpdate
    rollingUpdate:
      maxSurge: 5
      maxUnavailable: 10%
[...]
```

desired number of pods set to 10
max unavailable set to 30 percent

Settings Chat

Rolling Updates | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

An example of a rolling update strategy

Deployment

Desired Pods = 10 Pods
Max unavailable = 10% of desired Pods
Max surge = 5 Pods

Total Pods = 10

Old ReplicaSet
Number of Pods = 10

New ReplicaSet
Number of Pods = 5

the old replica set has 10 pods the deployment will begin by creating

Chat

Rolling Updates | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

An example of a rolling update strategy

Deployment

Desired Pods = 10 Pods
Max unavailable = 10% of desired Pods
Max surge = 5 Pods

Total Pods = 15 (Max)

Old ReplicaSet
Number of Pods = 10 - 6 = 4

New ReplicaSet
Number of Pods = 5

seven pods
a total of eight pods can be removed

Chat

Rolling Updates | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

An example of a rolling update strategy

Deployment

Desired Pods = 10 Pods
Max unavailable = 10% of desired Pods
Max surge = 5 Pods

Total Pods = 14

Old ReplicaSet
Number of Pods = $10 - 6 = 4$

New ReplicaSet
Number of Pods = $5 + 5 = 10$

this creates a total of 10 pods in a new replica

Chat

Rolling Updates | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49534

Rolling Updates

Mark as Completed

Video Rolling Updates

Video Blue-Green Deployments

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

An example of a rolling update strategy

Deployment

Desired Pods = 10 Pods
Max unavailable = 10% of desired Pods
Max surge = 5 Pods

Total Pods = 10

Old ReplicaSet
Number of Pods = $4 - 4 = 0$

New ReplicaSet
Number of Pods = $5 + 5 = 10$

finally the remaining two pods in the old set are deleted

Chat

Windows taskbar: 5:44 PM, 1/1/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/video/49535

Blue-Green Deployments

Video player controls: Home, Back, Forward, Stop, Play/Pause, Volume, Mute, Subtitles, Closed Captions, Full Screen, InPrivate, ...

Video title: Applying a Recreate strategy

Video content: A man in a grey blazer and black shirt is speaking. A code snippet is displayed in a box:

```
[...]
kind: deployment
spec:
  replicas: 10
  strategy:
    type: Recreate
[...]
```

Video player controls: Chat, Minimize, Maximize, Close.

Left sidebar: Blue-Green Deployments, Canary Deployments, Quiz: Updating deployments, Video: Managing Deployments, Video: Lab Intro, Hands-On Lab: Creating Google Kubernetes Engine Deployments, Video: Jobs and CronJobs, Video: Parallel Jobs.

Windows taskbar: 5:45 PM, 1/1/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/video/49535

Blue-Green Deployments

Video player controls: Home, Back, Forward, Stop, Play/Pause, Volume, Mute, Subtitles, Closed Captions, Full Screen, InPrivate, ...

Video title: Service is a load balancing front end for Pods

Video content: A man in a grey blazer and black shirt is speaking. A diagram is displayed in a box:

```
graph TD; Cluster[Cluster] --- FrontendPod[Frontend Pod]; FrontendPod --- Service[Service]; Service --- BackendPod1[Backend Pod]; Service --- BackendPod2[Backend Pod]; Service --- BackendPod3[Backend Pod];
```

Video player controls: Chat, Minimize, Maximize, Close.

Left sidebar: Blue-Green Deployments, Canary Deployments, Quiz: Updating deployments, Video: Managing Deployments, Video: Lab Intro, Hands-On Lab: Creating Google Kubernetes Engine Deployments, Video: Jobs and CronJobs, Video: Parallel Jobs.

Blue-Green Deployments | 5:45 PM 1/1/2021 InPrivate

Blue/Green deployment strategy

Client

= version: v1

Service

= version: v2

Deployment (my-app-v1)

Deployment (my-app-v2)

ReplicaSet

Pod

Pod

Pod

Pod

Pod

Pod

Chat

Blue-Green Deployments | 5:46 PM 1/1/2021 InPrivate

Applying a blue/green deployment strategy

```
[...]
kind: Service
spec:
  selector:
    app: my-app
    version: v1
[...]
```

```
$ kubectl apply -f my-app-v2.yaml
```

```
[...]
kind: Service
spec:
  selector:
    app: my-app
    version: v2
[...]
```

```
$ kubectl patch service my-app-service -p \
'{\"spec\": {\"selector\": {\"version\": \"v2\"}}}'
```

Chat

Canary Deployments | Qwik 5:50 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49536

Canary Deployments

Mark as Completed

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Canary deployments

Canary deployment

Client → Service → Deployment (my-app-v1) → Pod, Pod, Pod

Deployment (my-app-v2) → Pod, Pod, Pod

in this example 100% of the application traffic

0:34 / 3:46 Chat

Canary Deployments | Qwik 5:51 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49536

Canary Deployments

Mark as Completed

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Canary deployment

Client → Service → Deployment (my-app-v1) → Pod, Pod, Pod

Deployment (my-app-v2) → Pod, Pod, Pod

dash v1 when the canary deployment starts

0:34 / 3:46 Chat

Canary Deployments | Qwik 5:52 PM 1/1/2021 InPrivate

Canary Deployments

Mark as Completed

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Applying a canary deployment

```
[...] kind: Service spec: selector: app: my-app [...]
```

```
$ kubectl apply -f my-app-v2.yaml
```

```
$ kubectl scale deploy/my-app-v2 --replicas=10
```

```
$ kubectl delete -f my-app-v1.yaml
```

Scaled down and eventually deleted with the canary update strategy a subset of users

Chat

Canary Deployments | Qwik 5:52 PM 1/1/2021 InPrivate

Canary Deployments

Mark as Completed

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Session affinity ensures that all client requests are sent to the same Pod

Client

Service

Deployment (my-app-v1)

ReplicaSet

Pod

Deployment (my-app-v2)

ReplicaSet

Pod

connect to any pod deployment this potentially

Chat

Canary Deployments | Qwik 5:54 PM 1/1/2021 InPrivate

Canary Deployments

Mark as Completed

Video Canary Deployments

Quiz Quiz: Updating deployments

Video Managing Deployments

Video Lab Intro

Hands-On Lab Creating Google Kubernetes Engine Deployments

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Rolling back a Deployment

```
$ kubectl rollout undo deployment [DEPLOYMENT_NAME]
```

```
$ kubectl rollout undo deployment [DEPLOYMENT_NAME] --to-revision=2
```

```
$ kubectl rollout history deployment [DEPLOYMENT_NAME] --revision=2
```

Clean up Policy:

- Default: 10 Revision
- Change: .spec.revisionHistoryLimit

creation dates by default the details of the 10

Chat

Quiz: Updating deployments | Qwik 6:00 PM 1/1/2021 InPrivate

Quiz: Updating deployments

Congratulations! You passed this assessment.

✓ 1. You want to have two versions of your application in production, but be able to switch all traffic between them. This is an example of which deployment strategy?

Canary deployment

Blue-green deployment

That's correct!

Rolling updates

✓ 2. In a rolling update strategy, you can define the "max unavailable" parameter as a percentage. A percentage of what?

The total number of Pods in the cluster.

Chat

Quiz: Updating deployments | https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49537

Quiz: Updating deployments

Quiz
Quiz: Updating deployments

Video
Managing Deployments

Video
Lab Intro

Hands-On Lab
Creating Google Kubernetes Engine Deployments

Video
Jobs and CronJobs

Video
Parallel Jobs

Video
CronJobs

Quiz
Quiz: Jobs

That is correct.

✓ 2. In a rolling update strategy, you can define the "max unavailable" parameter as a percentage. A percentage of what?

The total number of Pods in the cluster.

The total number of Pods in the new ReplicaSet.

The total number of Pods across all ReplicaSets.

✓ 3. You want to have two versions of your application in production, but be able to a small percentage of traffic to the newer version as a gradual test. This is an example of which deployment strategy?

Blue-green deployment

Rolling updates

Canary deployment

Chat

Managing Deployments | https://googlecourses.qwiklabs.com/course_sessions/96271/video/49538

Managing Deployments

Mark as Completed

Video
Managing Deployments

Video
Lab Intro

Hands-On Lab
Creating Google Kubernetes Engine Deployments

Video
Jobs and CronJobs

Video
Parallel Jobs

Video
CronJobs

Quiz
Quiz: Jobs

Video
Lab Intro

Monitoring a Deployment

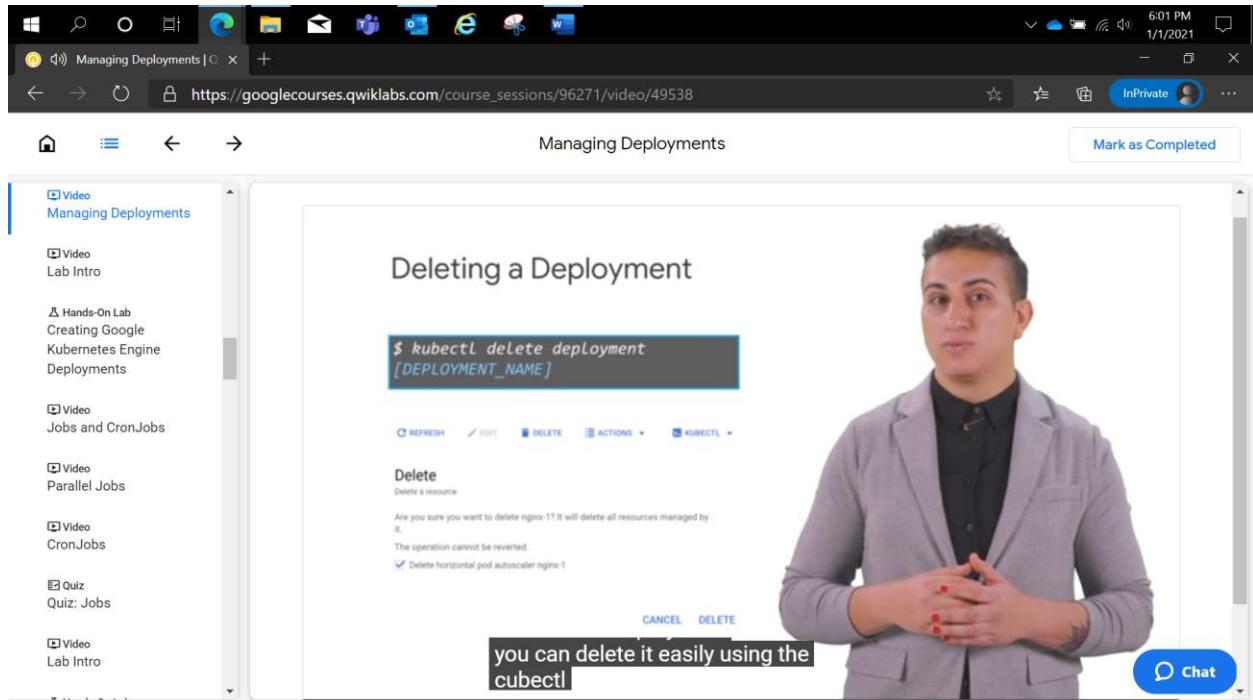
```
$ kubectl rollout pause deployment [DEPLOYMENT_NAME]
```

```
$ kubectl rollout resume deployment [DEPLOYMENT_NAME]
```

```
$ kubectl rollout status deployment [DEPLOYMENT_NAME]
```

you can also monitor the rollout status by using

Chat



Creating Google Kubernetes Engine Deployments

1 hour Free

Rate Lab

Overview

In this lab, you explore the basics of using deployment manifests. Manifests are files that contain configurations required for a deployment that can be used across different Pods. Manifests are easy to change.

Objectives

In this lab, you learn how to perform the following tasks:

- Create deployment manifests, deploy to cluster, and verify Pod rescheduling as nodes are disabled
- Trigger manual scaling up and down of Pods in deployments
- Trigger deployment rollout (rolling update to new version) and rollbacks
- Perform a Canary deployment

Task 0. Lab Setup

Access Qwiklabs

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an **incognito window**.

02:00:00

2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

START LAB

3. When ready, click **START LAB**.
4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

Open Google Console

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more.](#)

Username

google2876526_student@qwiklabs.n



Password

TG959yrKDX



GCP Project ID

qwiklabs-gcp-0855e773352d3560



[New to labs? View our introductory video!](#)

5. Click **Open Google Console**.
6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or **incur charges**.

7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it.

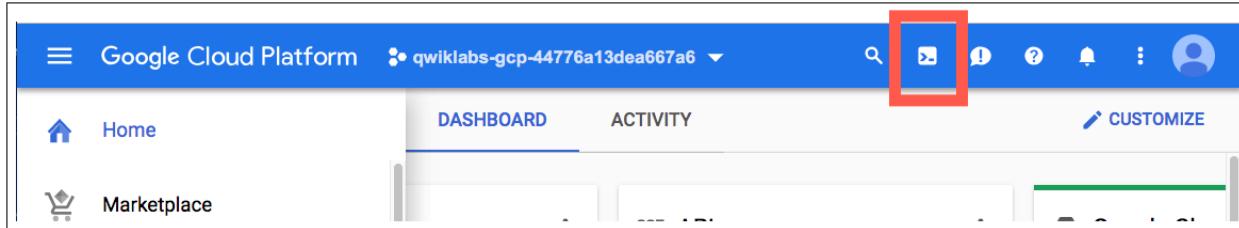
This clears your work and removes the project.

After you complete the initial sign-in steps, the project dashboard appears.

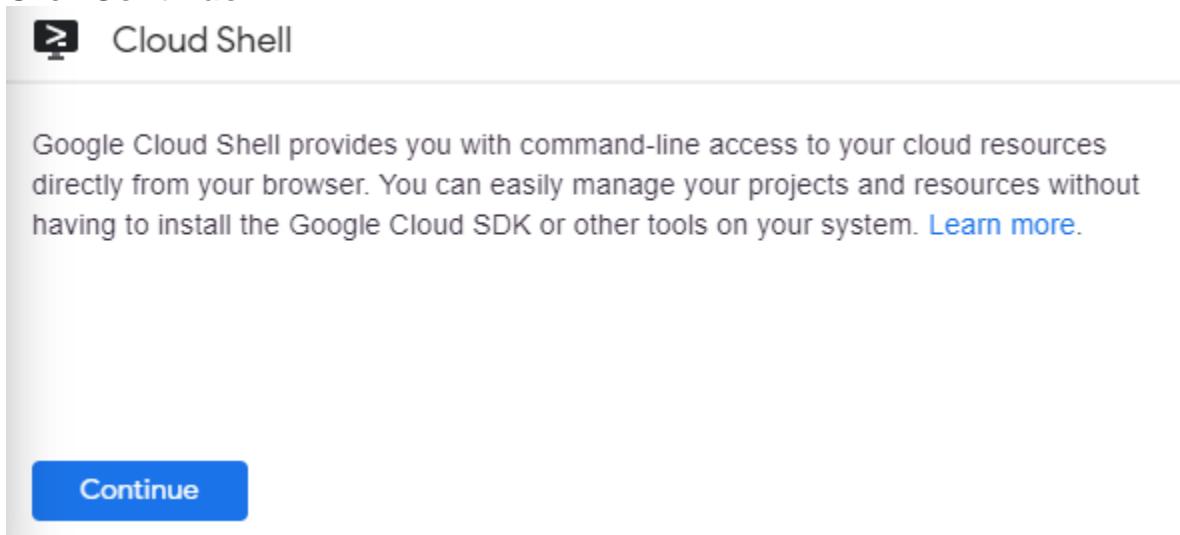
Activate Google Cloud Shell

Google Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Google Cloud Shell provides command-line access to your GCP resources.

1. In GCP console, on the top right toolbar, click the Open Cloud Shell button.



2. Click **Continue**.



It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your *PROJECT_ID*. For example:

```
Welcome to Cloud Shell! Type "help" to get started.  
Your Cloud Platform project in this session is set to qwiklabs-gcp-44776a13dea667a6.  
Use "gcloud config set project [PROJECT_ID]" to change to a different project.  
google1623327_student@cloudshell:~ (qwiklabs-gcp-44776a13dea667a6)$
```

gcloud is the command-line tool for Google Cloud Platform. It comes pre-installed on Cloud Shell and supports tab-completion.

You can list the active account name with this command:

```
gcloud auth list  
content_copy
```

Output:

```
Credentialed accounts:  
- <myaccount>@<mydomain>.com (active) content_copy
```

Example output:

```
Credentialed accounts:  
- google1623327 student@qwiklabs.netcontent_copy
```

You can list the project ID with this command:

```
gcloud config list project
```

```
content_copy
```

Output:

```
[core]
project = <project_ID>content_copy
```

Example output:

```
[core]
project = qwiklabs-gcp-44776a13dea667a6content_copy
```

Full documentation of **gcloud** is available on [Google Cloud gcloud Overview](#).

Task 1. Create deployment manifests and deploy to the cluster

In this task, you create a deployment manifest for a Pod inside the cluster.

Connect to the lab GKE cluster

1. In Cloud Shell, type the following command to set the environment variable for the zone and cluster name.

```
export my_zone=us-central1-a
export my_cluster=standard-cluster-1
content_copy
```

2. Configure kubectl tab completion in Cloud Shell.

```
source <(kubectl completion bash)
content_copy
```

3. In Cloud Shell, configure access to your cluster for the kubectl command-line tool, using the following command:

```
gcloud container clusters get-credentials $my_cluster --zone $my_zone
content_copy
```

4. In Cloud Shell enter the following command to clone the repository to the lab Cloud Shell.

```
git clone https://github.com/GoogleCloudPlatform/training-data-analyst  
content_copy
```

5. Create a soft link as a shortcut to the working directory.

```
ln -s ~/training-data-analyst/courses/ak8s/v1.1 ~/ak8s  
content_copy
```

6. Change to the directory that contains the sample files for this lab.

```
cd ~/ak8s/Deployments/  
content_copy
```

Create a deployment manifest

You will create a deployment using a sample deployment manifest called `nginx-deployment.yaml` that has been provided for you. This deployment is configured to run three Pod replicas with a single nginx container in each Pod listening on TCP port 80.

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-deployment
  labels:
    app: nginx
spec:
  replicas: 3
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
        - name: nginx
          image: nginx:1.7.9
          ports:
            - containerPort: 80
content_copy
```

1. To deploy your manifest, execute the following command:

```
kubectl apply -f ./nginx-deployment.yaml
```

`content_copy`

Click *Check my progress* to verify the objective.

Create and deploy manifest nginx deployment

Check my progress

2. To view a list of deployments, execute the following command:

```
kubectl get deployments  
content_copy
```

The output should look like this example.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
nginx-deployment	0/3	3	0	3s

`content_copy`

3. Wait a few seconds, and repeat the command until the number listed for CURRENT deployments reported by the command matches the number of DESIRED deployments.

The final output should look like the example.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
nginx-deployment	3/3	3	3	42s

`content_copy`

Task 2. Manually scale up and down the number of Pods in deployments

Sometimes, you want to shut down a Pod instance. Other times, you want ten Pods running. In Kubernetes, you can scale a specific Pod to the desired number of instances. To shut them down, you scale to zero.

In this task, you scale Pods up and down in the Google Cloud Console and Cloud Shell.

Scale Pods up and down in the console

1. Switch to the Google Cloud Console tab.



2. On the **Navigation menu** (), click **Kubernetes Engine > Workloads**.
3. Click **nginx-deployment** (your deployment) to open the Deployment details page.
4. At the top, click **ACTIONS > Scale**.

A screenshot of the Google Cloud Platform Kubernetes Engine Workloads page. The left sidebar shows 'Clusters' (selected), 'Workloads' (selected), 'Services', 'Applications', 'Configuration', and 'Storage'. The main area shows 'Deployment details' for 'nginx-deployment'. It includes a note to expose the deployment, three monitoring charts for CPU, Memory, and Disk usage over the last hour, and detailed deployment metadata: Cluster: standard-cluster-1, Namespace: default, Labels: app:nginx, Logs: Container logs, Audit logs, Replicas: 3 updated, 3 ready, 3 available, 0 unavailable, Pod specification: Revision 1, containers: nginx. Below this are sections for Active revisions (Revision 1, Name: nginx-deployment-75675f5897, Status: OK, Created on: Dec 16, 2018, 8:59:44 PM, Pods running/Pods total: 3/3) and Managed pods (Revision 1, Name: nginx-deployment-75675f5897-7bmix, Status: Running, Restarts: 0, Created on: Dec 16, 2018, 8:59:44 PM; Revision 1, Name: nginx-deployment-75675f5897-pfxqw, Status: Running, Restarts: 0, Created on: Dec 16, 2018, 8:59:44 PM). The top navigation bar has tabs for REFRESH, EDIT, DELETE, ACTIONS (with options: Autoscale, Expose, Rolling Update, Scale), and KUBECTL. A 'Expose' button is also present.

5. Type **1** and click **SCALE**.

This action scales down your cluster. You should see the Pod status being updated under **Managed Pods**. You might have to click **Refresh**.

Scale Pods up and down in the shell

1. Switch back to the Cloud Shell browser tab.
2. In the Cloud Shell, to view a list of Pods in the deployments, execute the following command:

```
kubectl get deployments  
content_copy
```

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
nginx-deployment	1/1	1	1	3m
content_copy				

3. To scale the Pod back up to three replicas, execute the following command:

```
kubectl scale --replicas=3 deployment nginx-deployment  
content_copy
```

4. To view a list of Pods in the deployments, execute the following command:

```
kubectl get deployments  
content_copy
```

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
nginx-deployment	3/3	3	3	4m
content_copy				

Task 3. Trigger a deployment rollout and a deployment rollback

A deployment's rollout is triggered if and only if the deployment's Pod template (that is, `.spec.template`) is changed, for example, if the labels or container

images of the template are updated. Other updates, such as scaling the deployment, do not trigger a rollout.

In this task, you trigger deployment rollout, and then you trigger deployment rollback.

Trigger a deployment rollout

1. To update the version of nginx in the deployment, execute the following command:

```
kubectl set image deployment.v1.apps/nginx-deployment nginx=nginx:1.9.1 --  
record  
content_copy
```

This updates the container image in your Deployment to nginx v1.9.1.

Click *Check my progress* to verify the objective.

Update version of nginx in the deployment

Check my progress

2. To view the rollout status, execute the following command:

```
kubectl rollout status deployment.v1.apps/nginx-deployment  
content_copy
```

The output should look like the example.

Output (do not copy)

```
Waiting for rollout to finish: 1 out of 3 new replicas updated...  
Waiting for rollout to finish: 1 out of 3 new replicas updated...  
Waiting for rollout to finish: 1 out of 3 new replicas updated...  
Waiting for rollout to finish: 2 out of 3 new replicas updated...  
Waiting for rollout to finish: 2 out of 3 new replicas updated...  
Waiting for rollout to finish: 2 out of 3 new replicas updated...  
Waiting for rollout to finish: 1 old replicas pending termination...  
Waiting for rollout to finish: 1 old replicas pending termination...  
deployment "nginx-deployment" successfully rolled out  
content_copy
```

3. To verify the change, get the list of deployments.

```
kubectl get deployments  
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
nginx-deployment	3/3	3	3	6m
content_copy				

4. View the rollout history of the deployment.

```
kubectl rollout history deployment nginx-deployment  
content_copy
```

The output should look like the example. Your output might not be an exact match.

Output (do not copy)

```
deployments "nginx-deployment"  
REVISION  CHANGE-CAUSE  
1          <none>  
2          kubectl set image deployment.v1.apps/nginx-deployment  
nginx=nginx:1.9.1 --record=true  
content_copy
```

Trigger a deployment rollback

To roll back an object's rollout, you can use the `kubectl rollout undo` command.

1. To roll back to the previous version of the nginx deployment, execute the following command:

```
kubectl rollout undo deployments nginx-deployment  
content_copy
```

2. View the updated rollout history of the deployment.

```
kubectl rollout history deployment nginx-deployment  
content_copy
```

The output should look like the example. Your output might not be an exact match.

Output (do not copy)

```
deployments "nginx-deployment"
REVISION CHANGE-CAUSE
2          kubectl set image deployment.v1.apps/nginx-deployment
nginx=nginx:1.9.1 --record=true
3          <none>
content_copy
```

3. View the details of the latest deployment revision

```
kubectl rollout history deployment/nginx-deployment --revision=3
content_copy
```

The output should look like the example. Your output might not be an exact match but it will show that the current revision has rolled back to nginx:1.7.9.

Output (do not copy)

```
deployments "nginx-deployment" with revision #3
Pod Template:
  Labels:      app=nginx
                pod-template-hash=3123191453
  Containers:
    nginx:
      Image:      nginx:1.7.9
      Port:       80/TCP
      Host Port:  0/TCP
      Environment:   <none>
      Mounts:     <none>
      Volumes:    <none>
content_copy
```

Task 4. Define the service type in the manifest

In this task, you create and verify a service that controls inbound traffic to an application. Services can be configured as ClusterIP, NodePort or LoadBalancer types. In this lab you configure a LoadBalancer.

Define service types in the manifest

A manifest file called `service-nginx.yaml` that deploys a LoadBalancer service type has been provided for you. This service is configured to distribute inbound traffic on TCP port 60000 to port 80 on any containers that have the label `app: nginx`.

```
apiVersion: v1
kind: Service
metadata:
  name: nginx
spec:
  type: LoadBalancer
  selector:
    app: nginx
  ports:
  - protocol: TCP
    port: 60000
    targetPort: 80
content_copy
```

1. In the Cloud Shell, to deploy your manifest, execute the following command:

```
kubectl apply -f ./service-nginx.yaml
content_copy
```

This manifest defines a service and applies it to Pods that correspond to the selector. In this case, the manifest is applied to the `nginx` container that you deployed in task 1. This service also applies to any other Pods with the `app: nginx` label, including any that are created after the service.

Click *Check my progress* to verify the objective.

Deploy manifest file that deploys LoadBalancer service type

[Check my progress](#)

Verify the LoadBalancer creation

1. To view the details of the `nginx` service, execute the following command:

```
kubectl get service nginx
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	CLUSTER_IP	EXTERNAL_IP	PORT(S)	SELECTOR	AGE
nginx	10.X.X.X	X.X.X.X	60000/TCP	run=nginx	1m
content_copy					

- When the external IP appears, open `http://[EXTERNAL_IP]:60000/` in a new browser tab to see the server being served through network load balancing.

It may take a few seconds before the **ExternalIP** field is populated for your service. This is normal. Just re-run the `kubectl get services nginx` command every few seconds until the field is populated.

Task 5. Perform a canary deployment

A canary deployment is a separate deployment used to test a new version of your application. A single service targets both the canary and the normal deployments. And it can direct a subset of users to the canary version to mitigate the risk of new releases. The manifest file `nginx-canary.yaml` that is provided for you deploys a single pod running a newer version of nginx than your main deployment. In this task, you create a canary deployment using this new deployment file.

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-canary
  labels:
    app: nginx
spec:
  replicas: 1
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
        track: canary
        Version: 1.9.1
```

```
spec:  
  containers:  
    - name: nginx  
      image: nginx:1.9.1  
      ports:  
        - containerPort: 80  
content_copy
```

The manifest for the nginx Service you deployed in the previous task uses a label selector to target the Pods with the `app: nginx` label. Both the normal deployment and this new canary deployment have the `app: nginx` label. Inbound connections will be distributed by the service to both the normal and canary deployment Pods. The canary deployment has fewer replicas (Pods) than the normal deployment, and thus it is available to fewer users than the normal deployment.

1. Create the canary deployment based on the configuration file.

```
kubectl apply -f nginx-canary.yaml  
content_copy
```

Click *Check my progress* to verify the objective.

Create a Canary Deployment

Check my progress

2. When the deployment is complete, verify that both the nginx and the nginx-canary deployments are present.

```
kubectl get deployments  
content_copy
```

3. Switch back to the browser tab that is connected to the external LoadBalancer service ip and refresh the page. You should continue to see the standard "Welcome to nginx" page.
4. Switch back to the Cloud Shell and scale down the primary deployment to 0 replicas.

```
kubectl scale --replicas=0 deployment nginx-deployment  
content_copy
```

5. Verify that the only running replica is now the Canary deployment:

```
kubectl get deployments  
content_copy
```

6. Switch back to the browser tab that is connected to the external LoadBalancer service ip and refresh the page. You should continue to see the standard "Welcome to nginx" page showing that the Service is automatically balancing traffic to the canary deployment.

Note: Session affinity

The Service configuration used in the lab does not ensure that all requests from a single client will always connect to the same Pod. Each request is treated separately and can connect to either the normal nginx deployment or to the nginx-canary deployment. This potential to switch between different versions may cause problems if there are significant changes in functionality in the canary release. To prevent this you can set the `sessionAffinity` field to `ClientIP` in the specification of the service if you need a client's first request to determine which Pod will be used for all subsequent connections.

For example:

```
apiVersion: v1
kind: Service
metadata:
  name: nginx
spec:
  type: LoadBalancer
  sessionAffinity: ClientIP
  selector:
    app: nginx
  ports:
  - protocol: TCP
    port: 60000
    targetPort: 80
```

The screenshot shows a Microsoft Edge browser window. The address bar displays the URL https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541. The main content area shows a video player with the title "Jobs and CronJobs". To the left of the video player is a sidebar containing a navigation menu with various video and quiz items. The central content area contains a diagram titled "A scenario where Job provides the solution". The diagram illustrates a flow from a user to a web server, then to a job pod labeled "Job: Transcode video files", and finally to a cluster containing a master node and two worker nodes. A red X marks the "Job Pod Container" on one of the worker nodes, indicating it has failed. A callout box states "and the job controller monitors the pod if a node failure occurs".

Jobs and CronJobs | Qwiklabs 6:38 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541

Jobs and CronJobs

Mark as Completed

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Deads

A scenario where Job provides the solution

The job was not able to run on the first node because it was busy. The job controller rescheduled the job to run on the second node.

Chat

Jobs and CronJobs | Qwiklabs 6:39 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541

Jobs and CronJobs

Mark as Completed

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Deads

A scenario where Job provides the solution

on a different node the job controller continues

Copy link

Chat

Jobs and CronJobs | Qwiklabs 6:40 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541

Jobs and CronJobs

Mark as Completed

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Details

A Job computing π to 2000 places

```
apiVersion: batch/v1
kind: Job
metadata:
  name: pi
spec:
  template:
    spec:
      containers:
        - name: pi
          image: perl
          command: ["perl", "-Mbignum=bpi", "-wle", "print bpi(2000)"]
          restartPolicy: Never
          backoffLimit: 4
```

this means that if a container in a pod fails for any reason, the job will fail.

Chat

Jobs and CronJobs | Qwiklabs 6:41 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541

Jobs and CronJobs

Mark as Completed

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Details

A Job computing π to 2000 places

```
apiVersion: batch/v1
kind: Job
metadata:
  name: pi
spec:
  template:
    spec:
      containers:
        - name: pi
          image: perl
          command: ["perl", "-Mbignum=bpi", "-wle", "print bpi(2000)"]
          restartPolicy: Never
          backoffLimit: 4
```

\$ kubectl apply -f [JOB_FILE]

\$ kubectl run pi --image perl --restart Never -- perl -Mbignum bpi -wle 'print bpi(2000)'

apply command

alternatively a job can be created using

Chat

Jobs and CronJobs | Qwiklabs 6:42 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49541

Jobs and CronJobs

Mark as Completed

Video Jobs and CronJobs

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

The role of a non-parallel Job

```
graph TD; A[.yaml file] --> B[Job object]; B --> C[Job controller]; C --> D[Pod or task]; D --> E["Job finished successfully"]
```

Pod or task
the pod is created as usual and the job

Chat

Parallel Jobs | Qwiklabs 6:48 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49542

Parallel Jobs

Mark as Completed

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Parallel Job with fixed completion count

```
apiVersion: batch/v1
kind: Job
metadata:
  name: my-app-job
spec:
  completions: 3
  parallelism: 2
  template:
    spec:
      [...]
```

```
graph TD; A[.yaml file] --> B[Job object]; B --> C[Job controller]
```

completions count is reached
in this example the controller will
launch

Chat

Parallel Jobs | Qwiklabs 6:49 PM 1/1/2021 InPrivate

Parallel Job with a worker queue

apiVersion: batch/v1
kind: Job
metadata:
 name: my-app-job
spec:
 parallelism: 3
 template:
 spec:
 [...]

```
graph TD; A[.yaml file] --> B[Job object]; B --> C[Job controller]
```

queue let's look at that next
in a worker queue parallel job

Chat

Parallel Jobs | Qwiklabs 6:50 PM 1/1/2021 InPrivate

Inspecting a Job

```
$ kubectl describe job [JOB_NAME]  
$ kubectl get pod -l [job-name=my-app-job]
```

describe command
the pods can be filtered using the cube

Chat

Parallel Jobs | Qwiklabs 6:50 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49542

Parallel Jobs

Mark as Completed

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz Cluster scaling

Scaling a Job

```
$ kubectl scale job [JOB_NAME] --replicas [VALUE]
```

Scale

Scale a workload to a new size.

Replicas

jobs can be scaled either from a command line or the gcp console

Chat



Parallel Jobs | Qwiklabs 6:51 PM 1/1/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49542

Parallel Jobs

Mark as Completed

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz Cluster scaling

Failing a Job

```
apiVersion: batch/v1
kind: Job
metadata:
  name: my-app-job
spec:
  backoffLimit: 4
  activeDeadlineSeconds: 300
  template:
  [...]
```

pods can fail all the time one way to limit pod failure

Chat



Parallel Jobs | Qwiklabs 6:51 PM 1/1/2021 InPrivate

Parallel Jobs

Mark as Completed

Video Parallel Jobs

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz Quiz: Cluster scaling

Retaining Job Pods

```
$ kubectl delete -f [JOB_FILE]
```

```
$ kubectl delete job [JOB_NAME]
```

```
$ kubectl delete job [JOB_NAME]
--cascade false
```

are also deleted if you want to retain



Chat

CronJobs | Qwiklabs 6:56 PM 1/1/2021 InPrivate

CronJobs

Mark as Completed

Video CronJobs

Quiz Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz Quiz: Cluster scaling

Video Controlling Pod Disruption

Setting up a cron schedule under the CronJob spec

```
apiVersion: batch/v1
kind: CronJob
metadata:
  name: my-app-job
spec:
  schedule: "*/1 * * * *"
  jobTemplate:
    spec:
      template:
        spec:
[ ... ]
```

scheduling a process the schedule field except a time in the unix/linux



Chat

6:56 PM
1/1/2021

InPrivate

Mark as Completed

CronJobs

Setting up a cron schedule under the CronJob spec

```
apiVersion: batch/v1
kind: CronJob
metadata:
  name: my-app-job
spec:
  schedule: "*/1 * * * *"
  startingDeadlineSeconds: 3600
  jobTemplate:
    spec:
      template:
        spec:
          [...]
```

about jobs if the job defined in this manifest is scheduled

Chat

Video CronJobs

Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz: Cluster scaling

Video Controlling Pod Placement

6:58 PM
1/1/2021

InPrivate

Mark as Completed

CronJobs

Setting up a cron schedule under the CronJob spec

```
apiVersion: batch/v1
kind: CronJob
metadata:
  name: my-app-job
spec:
  schedule: "*/1 * * * *"
  startingDeadlineSeconds: 3600
  concurrencyPolicy: Forbid
  suspend: True
  successfulJobsHistoryLimit: 3
  failedJobsHistoryLimit: 1
  jobTemplate:
    spec:
      [...]
```

considered or directed you can stop execution of individual jobs by a cron job by

Chat

Video CronJobs

Quiz: Jobs

Video Lab Intro

Hands-On Lab Deploying Jobs on Google Kubernetes Engine

Video Cluster Scaling

Video Downscaling

Video Node Pools

Quiz: Cluster scaling

Video Controlling Pod Placement

Windows taskbar: 6:58 PM, 1/1/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/video/49543

CronJobs

Mark as Completed

Video player controls: Home, List, Previous, Next, Stop, Volume, Mute, Subtitles, Fullscreen, Closed Captions.

Video content: A man in a grey blazer and black shirt is speaking. Text overlay: "CronJobs can be managed using kubectl".

Left sidebar:

- Video CronJobs
- Quiz: Jobs
- Video Lab Intro
- Hands-On Lab Deploying Jobs on Google Kubernetes Engine
- Video Cluster Scaling
- Video Downscaling
- Video Node Pools
- Quiz: Cluster scaling
- Video Controlling Pod Placement

Code snippets:

- Create a CronJob: `$ kubectl apply -f [FILE]`
- Inspect a CronJob: `$ kubectl describe cronjob [NAME]`
- Delete a CronJob: `$ kubectl delete cronjob [NAME]`

Text at bottom: "failed jobs history limit cron jobs operate in the same manner as the job"

Chat button: Chat

Windows taskbar: 6:59 PM, 1/1/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49544

Quiz: Jobs

Your score: 100% Passing score: 50% Retake

Congratulations! You passed this assessment.

Left sidebar:

- Quiz: Jobs
- Video Lab Intro
- Hands-On Lab Deploying Jobs on Google Kubernetes Engine
- Video Cluster Scaling
- Video Downscaling
- Video Node Pools
- Quiz: Cluster scaling
- Video Controlling Pod Placement
- Video Affinity and Anti-Affinity

Question 1:

✓ 1. What happens if a node fails while a Job is executing on that node?

Kubernetes will wait for the node to return to service and then restart the Job.

Kubernetes will abandon the Job.

Kubernetes will restart the Job on a node that is still running.

That is correct.

Question 2:

✓ 2. Suppose you have a Job in which each Pod performs work drawn from a work queue. How should this Job's manifest be configured?

Specify a spec.completions value and leave the parallelism value

Chat button: Chat

The screenshot shows a Microsoft Edge browser window with the following details:

- Title Bar:** Shows the URL https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49544, the page title "Quiz: Jobs", and the date/time "1/1/2021 6:59 PM".
- Header:** "Quiz: Jobs"
- Left Sidebar:** A vertical list of course content:
 - Quiz: Jobs
 - Video Lab Intro
 - Hands-On Lab Deploying Jobs on Google Kubernetes Engine
 - Video Cluster Scaling
 - Video Downscaling
 - Video Node Pools
 - Quiz Quiz: Cluster scaling
 - Video Controlling Pod Placement
 - Video Affinity and Anti-Affinity
- Main Content Area:** Displays a question:

That is correct.

✓ 2. Suppose you have a Job in which each Pod performs work drawn from a work queue. How should this Job's manifest be configured?

 - Specify a spec.completions value and leave the parallelism value unset
 - Specify a WorkQueue object to let Kubernetes query the state of the queue.
 - Specify a parallelism value and leave spec.completions unset

That is correct.
- Bottom Right:** A blue "Chat" button.

Deploying Jobs on Google Kubernetes Engine

1 hourFree

Rate Lab

Overview

In this lab, you define and run Jobs and CronJobs.

In GKE, a Job is a controller object that represents a finite task. Jobs manage a task as it runs to completion, rather than managing an ongoing desired state such as the maintaining the total number of running Pods.

CronJobs perform finite, time-related tasks that run once or repeatedly at a time that you specify using Job objects to complete their tasks.

Objectives

In this lab, you learn how to perform the following tasks:

- Define, deploy and clean up a GKE Job
- Define, deploy and clean up a GKE CronJob

Task 0. Lab Setup

Access Qwiklabs

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an **incognito window**.

02:00:00

2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

START LAB

3. When ready, click **START LAB**.
4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

Open Google Console

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more.](#)

Username

google2876526_student@qwiklabs.n



Password

TG959yrKDX



GCP Project ID

qwiklabs-gcp-0855e773352d3560



[New to labs? View our introductory video!](#)

5. Click **Open Google Console**.
6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or **incur charges**.

7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it.

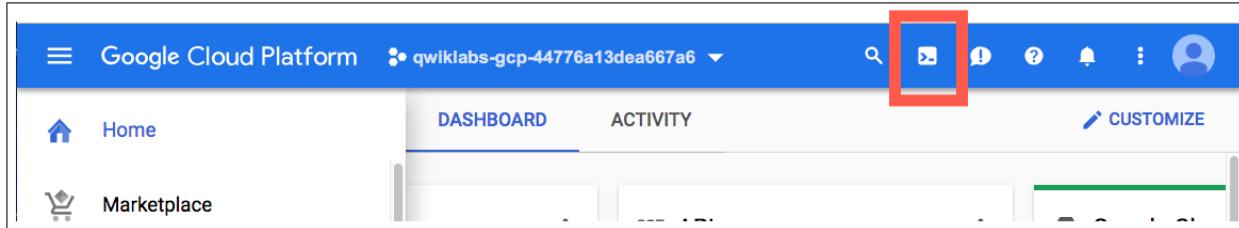
This clears your work and removes the project.

After you complete the initial sign-in steps, the project dashboard appears.

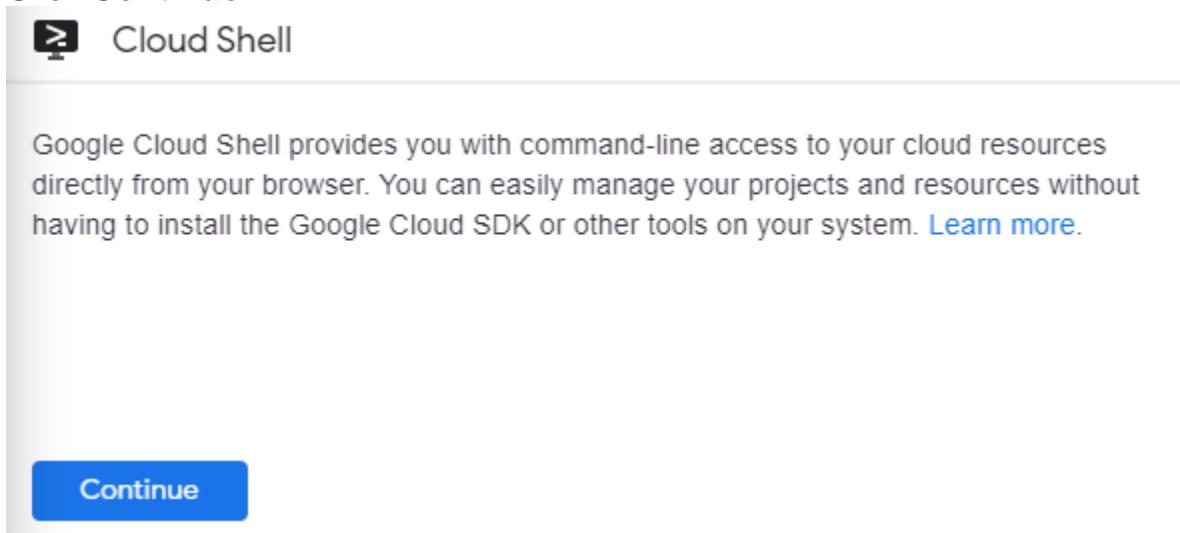
Activate Google Cloud Shell

Google Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Google Cloud Shell provides command-line access to your GCP resources.

1. In GCP console, on the top right toolbar, click the Open Cloud Shell button.



2. Click **Continue**.



It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your *PROJECT_ID*. For example:

A screenshot of a Cloud Shell terminal window. The title bar shows "...abs-gcp-44776a13dea667a6". The terminal output shows a welcome message and the active project ID: "Your Cloud Platform project in this session is set to **qwiklabs-gcp-44776a13dea667a6**". A red arrow points to the project ID text. The prompt at the bottom is "google1623327_student@cloudshell:~ (qwiklabs-gcp-44776a13dea667a6)\$".

gcloud is the command-line tool for Google Cloud Platform. It comes pre-installed on Cloud Shell and supports tab-completion.

You can list the active account name with this command:

```
gcloud auth list  
content_copy
```

Output:

```
Credentialed accounts:  
- <myaccount>@<mydomain>.com (active) content_copy
```

Example output:

```
Credentialed accounts:  
- google1623327 student@qwiklabs.netcontent_copy
```

You can list the project ID with this command:

```
gcloud config list project
```

```
content_copy
```

Output:

```
[core]
project = <project_ID>content_copy
```

Example output:

```
[core]
project = qwiklabs-gcp-44776a13dea667a6content_copy
```

Full documentation of **gcloud** is available on [Google Cloud gcloud Overview](#).

Task 1. Define and deploy a Job manifest

In GKE, a Job is a controller object that represents a finite task.

In this task, you create a Job, inspect its status, and then remove it.

Connect to the lab Google Kubernetes Engine cluster

1. In Cloud Shell, type the following command to set the environment variable for the zone and cluster name.

```
export my_zone=us-central1-a
export my_cluster=standard-cluster-1
content_copy
```

2. Configure kubectl tab completion in Cloud Shell.

```
source <(kubectl completion bash)
content_copy
```

3. In Cloud Shell, configure access to your cluster for the kubectl command-line tool, using the following command:

```
gcloud container clusters get-credentials $my_cluster --zone $my_zone
```

```
content_copy
```

4. In Cloud Shell enter the following command to clone the repository to the lab Cloud Shell.

```
git clone https://github.com/GoogleCloudPlatform/training-data-analyst  
content_copy
```

5. Create a soft link as a shortcut to the working directory.

```
ln -s ~/training-data-analyst/courses/ak8s/v1.1 ~/ak8s  
content_copy
```

6. Change to the directory that contains the sample files for this lab.

```
cd ~/ak8s/Jobs_CronJobs  
content_copy
```

Create and run a Job

You will create a job using a sample deployment manifest called `example-job.yaml` that has been provided for you. This Job computes the value of Pi to 2,000 places and then prints the result.

```
apiVersion: batch/v1
kind: Job
metadata:
  # Unique key of the Job instance
  name: example-job
spec:
  template:
    metadata:
      name: example-job
    spec:
      containers:
        - name: pi
          image: perl
          command: ["perl"]
          args: ["-Mbignum=bpi", "-wle", "print bpi(2000)"]
          # Do not restart containers after they exit
          restartPolicy: Never
content_copy
```

1. To create a Job from this file, execute the following command:

```
kubectl apply -f example-job.yaml  
content_copy
```

Click *Check my progress* to verify the objective.

Create and run a Job

Check my progress

2. To check the status of this Job, execute the following command:

```
kubectl describe job example-job  
content_copy
```

You will see details of the job, including the Pod statuses indicating how many jobs are still running, how many completed successfully and how many failed.

Output (do not copy)

```
...  
Start Time: Thu, 20 Dec 2018 14:34:09 +0000  
Pods Statuses: 0 Running / 1 Succeeded / 0 Failed  
...  
content_copy
```

3. To view all Pod resources in your cluster, including Pods created by the Job which have completed, execute the following command:

```
kubectl get pods  
content_copy
```

Your Pod name might be different from the example output. Make a note of one of the Pod names.

Output (do not copy)

NAME	READY	STATUS	RESTARTS	AGE
example-job-sqljc	0/1	Completed	0	1m

```
content_copy
```

Clean up and delete the Job

When a Job completes, the Job stops creating Pods. The Job API object is not removed when it completes, which allows you to view its status. Pods created by the Job are not deleted, but they are terminated. Retention of the Pods allows you to view their logs and to interact with them.

1. To get a list of the Jobs in the cluster, execute the following command:

```
kubectl get jobs  
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	COMPLETIONS	DURATION	AGE
example-job	1/1	75s	2m5s
content_copy			

2. To retrieve the log file from the Pod that ran the Job execute the following command. You must replace [POD-NAME] with the node name you recorded in the last task

```
kubectl logs [POD-NAME]  
content_copy
```

The output will show that the job wrote the first two thousand digits of pi to the Pod log.

3. To delete the Job, execute the following command:

```
kubectl delete job example-job  
content_copy
```

If you try to query the logs again the command will fail as the Pod can no longer be found.

Task 2. Define and deploy a CronJob manifest

You can create CronJobs to perform finite, time-related tasks that run once or repeatedly at a time that you specify.

In this task, you create and run a CronJob, and then you clean up and delete the Job.

Create and run a CronJob

The CronJob manifest file `example-cronjob.yaml` has been provided for you. This CronJob deploys a new container every minute that prints the time, date and "Hello, World!".

```
apiVersion: batch/v1beta1
kind: CronJob
metadata:
  name: hello
spec:
  schedule: "*/1 * * * *"
  jobTemplate:
    spec:
      template:
        spec:
          containers:
            - name: hello
              image: busybox
              args:
                - /bin/sh
                - -c
                - date; echo "Hello, World!"
          restartPolicy: OnFailure
```

content_copy

Note

CronJobs use the required `schedule` field, which accepts a time in the Unix standard `crontab` format. All CronJob times are in UTC:

- The first value indicates the minute (between 0 and 59).
 - The second value indicates the hour (between 0 and 23).
 - The third value indicates the day of the month (between 1 and 31).
 - The fourth value indicates the month (between 1 and 12).
 - The fifth value indicates the day of the week (between 0 and 6).
- The `schedule` field also accepts * and ? as wildcard values. Combining / with ranges specifies that the task should repeat at a regular interval. In the example, `*/1 * * * *` indicates that the task should repeat every minute of every day of every month.

1. To create a Job from this file, execute the following command:

```
kubectl apply -f example-cronjob.yaml
```

content_copy

Click *Check my progress* to verify the objective.

Create and run a CronJob

Check my progress

- To get a list of the Jobs in the cluster, execute the following command:

```
kubectl get jobs  
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	COMPLETIONS	DURATION	AGE
hello-1545013620	1/1	2s	18s
content_copy			

- To check the status of this Job, execute the following command, where [job_name] is the name of your job:

```
kubectl describe job [job_name]  
content_copy
```

You will see details of the job, including the Pod statuses showing that one instance of this job was run.

Output (do not copy)

```
...  
Start Time: Thu, 20 Dec 2018 15:24:03 +0000  
Pods Statuses: 0 Running / 1 Succeeded / 0 Failed  
...  
...Created pod: hello-1545319920-twkhl  
content_copy
```

- Make a note of the name of the Pod that was used by this job.

- View the output of the Job by querying the logs for the Pod. Replace [POD-NAME] with the name of the Pod you recorded in the last step.

```
kubectl logs [POD-NAME]  
content_copy
```

This will display the output of the shell script configured in the CronJob:

Output (do not copy)

```
Thu Dec 20 15:31:16 WET 2018  
Hello,World!  
content_copy
```

- To view all job resources in your cluster, including all of the Pods created by the CronJob which have completed, execute the following command:

```
kubectl get jobs  
content_copy
```

Your job names might be different from the example output. By default Kubernetes sets the Job history limits so that only the last three successful and last failed job are retained so this list will only contain the most recent three of four jobs.

Output (do not copy)

NAME	COMPLETIONS	DURATION	AGE
hello-1545013680	1	1	2m
hello-1545013740	1	1	1m
hello-1545013800	1	1	42s

```
content_copy
```

Clean up and delete the Job

In order to stop the CronJob and clean up the Jobs associated with it you must delete the CronJob.

1. To delete all these jobs, execute the following command:

```
kubectl delete cronjob hello  
content_copy
```

2. To verify that the jobs were deleted, execute the following command:

```
kubectl get jobs  
content_copy
```

The output should look like the example.

Output (do not copy)

```
No resources found.  
content_copy
```

All the Jobs were removed

```
student_03_da03603be4b2@cloudshell:~/ak8s/Jobs_CronJobs (qwiklabs-gcp-03-363474b52250)$ cat example-job.yaml
apiVersion: batch/v1
kind: Job
metadata:
  # Unique key of the Job instance
  name: example-job
```

```
spec:  
template:  
  metadata:  
    name: example-job  
  spec:  
    containers:  
      - name: pi  
        image: perl  
        command: ["perl1"]  
        args: ["-Mbignum=bpi", "-wle", "print bpi(2000)"]  
      # Do not restart containers after they exit  
    restartPolicy: Never
```

The screenshot shows a Microsoft Edge browser window with the following details:

- Title Bar:** Cluster Scaling | Qwiklabs
- Address Bar:** https://googlecourses.qwiklabs.com/course_sessions/96271/video/49547
- Header:** Cluster Scaling (with a "Mark as Completed" button)
- Left Sidebar (Table of Contents):**
 - Video Parallel Jobs
 - Video CronJobs
 - Quiz Quiz: Jobs
 - Video Lab Intro
 - Hands-On Lab Deploying Jobs on Google Kubernetes Engine
 - Video Cluster Scaling (selected)
 - Video Downscaling
 - Video Node Pools
 - Quiz Quiz: Cluster scaling
- Content Area:**
 - Section Title:** Scaling down a cluster by using the gcloud command
 - Code Example:** `gcloud container clusters resize projectdemo --node-pool default-pool \ --size 6`
 - Text Callout:** nodes without pods aren't differentiated
resize will pick instances to remove at
 - Speaker:** A man in a grey blazer and black shirt is speaking.
 - Bottom Right:** Chat button

Cluster Scaling | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49547

Cluster Scaling

Mark as Completed

Scale up a cluster with autoscaling

Cluster

Master

Node

Pod

Container

other parts terminate free up capacity or additional nodes are added

Chat

Video Cluster Scaling

Parallel Jobs

CronJobs

Quiz: Jobs

Lab Intro

Hands-On Lab

Deploying Jobs on Google Kubernetes Engine

Downscaling

Node Pools

Quiz: Cluster scaling

A screenshot of a video player interface from Google Courses on Qwiklabs. The video title is "Scale up a cluster with autoscaling". On the left, there's a sidebar with a list of course modules: Parallel Jobs, CronJobs, Quiz: Jobs, Lab Intro, Hands-On Lab (disabled), Deploying Jobs on Google Kubernetes Engine, Cluster Scaling (selected), Downscaling, Node Pools, and Quiz: Cluster scaling. The main content area shows a woman in a grey blazer speaking. To her left is a diagram of a cluster with a master node and two blue node boxes. Each node box contains a pod with a container. A question mark icon is positioned next to the second node. A callout box at the bottom of the diagram says "other parts terminate free up capacity or additional nodes are added". A "Chat" button is in the bottom right corner.

Cluster Scaling | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49547

Cluster Scaling

Mark as Completed

Scale up a cluster with autoscaling

Cluster

Master

Node

Pod

Container

resources to become available and the pod is then scheduled

Chat

Pause (k)

Video Cluster Scaling

Parallel Jobs

CronJobs

Quiz: Jobs

Lab Intro

Hands-On Lab

Deploying Jobs on Google Kubernetes Engine

Downscaling

Node Pools

Quiz: Cluster scaling

A screenshot of the same video player interface, showing the next frame. The woman is still speaking. The diagram now shows the second node in green with a large white plus sign in the center. A callout box at the bottom of the diagram says "resources to become available and the pod is then scheduled". A "Pause (k)" button is visible at the bottom left of the diagram area.

Windows taskbar: 11:11 AM, 1/2/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/video/49548

Downscaling

Mark as Completed

Video player controls: Home, List, Back, Forward, Stop, Volume, Mute, Subtitles, Fullscreen, InPrivate, ...

Video title: Scale down a cluster with autoscaling

Video content:

- 1 There can be no scale-up events pending.
- 2 Can the node be deleted safely?

Text overlay: second it texts that the node can be deleted safely

Speaker icon: Chat

Video navigation sidebar:

- Video Downscaling
- Video Node Pools
- Quiz Quiz: Cluster scaling
- Video Controlling Pod Placement
- Video Affinity and Anti-Affinity
- Video Pod Placement Example
- Video Taints and Tolerations
- Quiz Quiz: Controlling pod placement
- Video Getting software into your

Windows taskbar: 11:12 AM, 1/2/2021, InPrivate

Address bar: https://googlecourses.qwiklabs.com/course_sessions/96271/video/49548

Downscaling

Mark as Completed

Video player controls: Home, List, Back, Forward, Stop, Volume, Mute, Subtitles, Fullscreen, InPrivate, ...

Video title: Pod conditions that prevent node deletion

Video content:

- Not run by a controller
- Has local storage
- Restricted by constraint rules

Text overlay: entered restricted by constraint rules that prevent

Speaker icon: Chat

Video navigation sidebar:

- Video Downscaling
- Video Node Pools
- Quiz Quiz: Cluster scaling
- Video Controlling Pod Placement
- Video Affinity and Anti-Affinity
- Video Pod Placement Example
- Video Taints and Tolerations
- Quiz Quiz: Controlling pod placement
- Video Getting software into your

Downscaling | Qwiklabs 11:12 AM 1/2/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49548

Downscaling

Mark as Completed

Video Downscaling

Video Node Pools

Quiz Quiz: Cluster scaling

Video Controlling Pod Placement

Video Affinity and Anti-Affinity

Video Pod Placement Example

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your

Pod conditions that prevent node deletion

- X** cluster-autoscaler.kubernetes.io/safe-to-evict is set to False
- X** Restrictive PodDisruptionBudget
- X** kubernetes.io/scale-down-disabled set to True

evicted or disrupted at a time at the node level if the node scale



Chat

Downscaling | Qwiklabs 11:14 AM 1/2/2021 InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49548

Downscaling

Mark as Completed

Video Downscaling

Video Node Pools

Quiz Quiz: Cluster scaling

Video Controlling Pod Placement

Video Affinity and Anti-Affinity

Video Pod Placement Example

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your

Best practices for working with autoscaled clusters

X

- Don't run Compute Engine autoscaling
- Don't manually resize a node using the gcloud command
- Don't manually modify a node

✓

- Specify correct resource requests for Pods
- Use PodDisruptionBudget to maintain the app's availability

controller can be safely terminated and relocated



Chat

Node Pools

Setting a node pool size



Node pool = 0
Cluster size ≠ 0

MAX = 1,000 nodes x 30 Pods

Increase quota limits to avoid disruption
with each running 30 pods
however standard DCP quota

Chat

Node Pools

Mark as Completed

gcloud commands for autoscaling

Create a cluster with autoscaling enabled	Enable autoscaling for an existing node pool
<pre>gcloud container clusters create [CLUSTER_NAME] --num-nodes 30 \ --enable-autoscaling --min-nodes 15 \ --max-nodes 50 [-z zone COMPUTE_ZONE]</pre>	<pre>gcloud container clusters update [CLUSTER_NAME] --enable-autoscaling \ --min-nodes 1 --max-nodes 10 --zone [COMPUTE_ZONE] --node-pool [POOL_NAME]</pre>
Add a node pool with autoscaling enabled	Disable autoscaling for an existing node pool
<pre>gcloud container node-pools create [POOL_NAME] --cluster [CLUSTER_NAME] \ --enable-autoscaling --min-nodes 15 \ --max-nodes 50 [-z zone COMPUTE_ZONE]</pre>	<pre>gcloud container clusters update [CLUSTER_NAME] --no-enable-autoscaling \ --node-pool [POOL_NAME] [-z zone [COMPUTE_ZONE]] --project [PROJECT_ID]</pre>

scaling for an existing node pool you can also disable auto scaling

Quiz: Cluster scaling | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49550

Quiz: Cluster scaling

Quiz: Cluster scaling

1. True or false: if you manually decrease the size of a node pool, any Pods on deleted nodes will be restarted on other nodes.

False

True

That is correct.

2. True or false: if autoscaling decreases the size of a node pool, any Pods on deleted nodes that aren't managed by a replication controller will be lost.

False

True

That is correct.

Chat

Quiz: Cluster scaling | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49550

Quiz: Cluster scaling

Quiz: Cluster scaling

False

True

That is correct.

3. True or false: if you increase the size of a node pool, existing Pods will be moved to newer nodes.

False

True

That is correct.

Chat

Controlling Pod Placement | InPrivate

Controlling Pod Placement

Controlled scheduling

when nodes are started the cubelet automatically

Chat

Mark as Completed

Quiz: Cluster scaling

Video: Controlling Pod Placement

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Tolerations

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Controlling Pod Placement | InPrivate

Controlling Pod Placement

Nodes must match all the labels present under the nodeSelector field

```
apiVersion: v1
kind: Pod
metadata:
  name: nginx
  labels:
    env: test
spec:
  containers:
    - name: nginx
      image: nginx
      imagePullPolicy: IfNotPresent
      nodeSelector:
        disktype: ssd
[...]
```

```
apiVersion: v1
kind: Node
metadata:
  name: node1
  labels:
    disktype: ssd
[...]
```

assigned for example the label kubernetes

Chat

Mark as Completed

Quiz: Cluster scaling

Video: Controlling Pod Placement

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Tolerations

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Controlling Pod Placement | InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49551

Controlling Pod Placement

Mark as Completed

Quiz: Cluster scaling

Video: Controlling Pod Placement

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Nodes must match all the labels present under the nodeSelector field

```
apiVersion: v1
kind: Pod
metadata:
  name: mysql
  labels:
    env: test
spec:
  containers:
    - name: mysql
      image: mysql
      imagePullPolicy: IfNotPresent
  nodeSelector:
    cloud.google.com/gke-nodepool=ssd
[...]
```

for example this pod will only run on a node that is labeled with a gke label

Chat

Affinity and Anti-Affinity | InPrivate

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49552

Affinity and Anti-Affinity

Mark as Completed

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Hands-On Lab: Configuring Pod Autoscaling and NodePools

Quiz: Deployments... jobs

Node affinity is conceptually similar to nodeSelector

```
apiVersion: v1
kind: Pod
metadata:
  name: with-node-affinity
spec:
  affinity:
    nodeAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        nodeSelectorTerms:
          - matchExpressions:
              - key: beta.kubernetes.io/instance-type
                operator: In
                values:
                  - n1-highmem-4
                  - n1-highmem-8
[...]
```

next let's look at no Definity in ante

Chat

Affinity and Anti-Affinity | [InPrivate](#)

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49552

Affinity and Anti-Affinity

Mark as Completed

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Hands-On Lab: Configuring Pod Autoscaling and NodePools

Defining the intensity of preference

```
apiVersion: v1
kind: Pod
metadata:
  name: with-node-affinity
spec:
  affinity:
    nodeAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
        - weight: 1
          preference:
            - matchExpressions:
                - key: accelerator-type
                  operator: In
                  values:
                    - gpu
                    - tpu
```

1 well it makes it a soft preference rather than the hard requirement is our

Chat

Affinity and Anti-Affinity | [InPrivate](#)

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49552

Affinity and Anti-Affinity

Mark as Completed

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Hands-On Lab: Configuring Pod Autoscaling and NodePools

Affinity and anti-affinity rules

```
apiVersion: v1
kind: Pod
metadata:
  name: with-node-affinity
spec:
  affinity:
    nodeAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        nodeSelectorTerms:
          - matchExpressions:
              - key: accelerator-type
                operator: In
                values:
                  - gpu
                  - tpu
```

TPU these could be any labels you want to use in this case the

Chat

Affinity and Anti-Affinity | [InPrivate](#)

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49552

Affinity and Anti-Affinity

Mark as Completed

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Hands-On Lab: Configuring Pod Autoscaling and NodePools

Inter-pod affinity and anti-affinity features are built on the node affinity concept

```
[...]
metadata:
  name: with-pod-affinity
spec:
  affinity:
    podAntiAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
        - weight: 100
          podAffinityTerm:
            labelSelector:
              - matchExpressions:
                  - key: app
                    operator: In
                  values:
                    - webserver
```

preference inter pod affinity in anti affinity features extend

Chat

Affinity and Anti-Affinity | [InPrivate](#)

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49552

Affinity and Anti-Affinity

Mark as Completed

Video: Affinity and Anti-Affinity

Video: Pod Placement Example

Video: Taints and Toleration

Quiz: Controlling pod placement

Video: Getting software into your cluster

Video: Lab Intro

Hands-On Lab: Configuring Pod Autoscaling and NodePools

Inter-pod affinity and anti-affinity features are built on the node affinity concept

```
[...]
metadata:
  name: with-pod-affinity
spec:
  affinity:
    podAntiAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
        - weight: 100
          podAffinityTerm:
            labelSelector:
              - matchExpressions:
                  - key: app
                    operator: In
                  values:
                    - webserver
```

topologyKey: failure-domain.beta.kubernetes.io/zone
topology keys you can also specify
affinity and non

Chat

Pod Placement Example | https://googlecourses.qwiklabs.com/course_sessions/96271/video/49553

Pod Placement Example

Combining inter-pod affinity and anti-affinity

Pod (#1)
(Label – app:webserver)
Requirement: No app:webserver on the same zone
Preference: Prefer app:cache on the same node

Pod (#2)
(Label – app:webserver)
Requirement: No app:webserver on the same zone
Preference: Prefer app:cache on the same node

Pod (#3)
(Label – app:cache)
Requirement: No app:cache on the same zone
Preference: Prefer app:webserver on the same node

Pod (#4)
(Label – app:cache)
Requirement: No app:cache on the same zone
Preference: Prefer app:webserver on the same node

Chat

Video
Pod Placement Example

Video
Taints and Tolerations

Quiz
Quiz: Controlling pod placement

Video
Getting software into your cluster

Video
Lab Intro

Hands-On Lab
Configuring Pod
Autoscaling and NodePools

Quiz
Quiz: Deployments, Jobs, and Scaling

Mark as Completed

Pod Placement Example | https://googlecourses.qwiklabs.com/course_sessions/96271/video/49553

Chapter 2.2.2 - Deployment, Jobs and Scaling

Combining inter-pod affinity and anti-affinity

Zone 1

Node

Pod (#1)
(Label – app:webserver)
Requirement: No app:webserver on the same zone
Preference: Prefer app:cache on the same node

Pod (#3)
(Label – app:cache)
Requirement: No app:cache on the same zone
Preference: Prefer app:webserver on the same node

Zone 2

Node

Pod (#2)
(Label – app:webserver)
Requirement: No app:webserver on the same zone
Preference: Prefer app:cache on the same node

Pod (#4)
(Label – app:cache)
Requirement: No app:cache on the same zone
Preference: Prefer app:webserver on the same node

Chat

Video
Pod Placement Example

Video
Taints and Tolerations

Quiz
Quiz: Controlling pod placement

Video
Getting software into your cluster

Video
Lab Intro

Hands-On Lab
Configuring Pod
Autoscaling and NodePools

Quiz
Quiz: Deployments, Jobs, and Scaling

Mark as Completed

Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Mark as Completed

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Node

\$ kubectl taint nodes node1 key=value:NoSchedule

Node

Tainted with:
Key=Value: NoSchedule
this tank with the no schedule effect limits all

Chat

Taints and Tolerations

Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Mark as Completed

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

You can apply multiple taints to a node

Node

Tainted with:
Key=Value: NoSchedule
Key2=Value2: NoExecute

in this case no pods can be scheduled and

Chat

Taints and Tolerations

11:43 AM 1/2/2021

Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Mark as Completed

InPrivate

Home Back Forward Taints and Tolerations

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Google Kubernetes

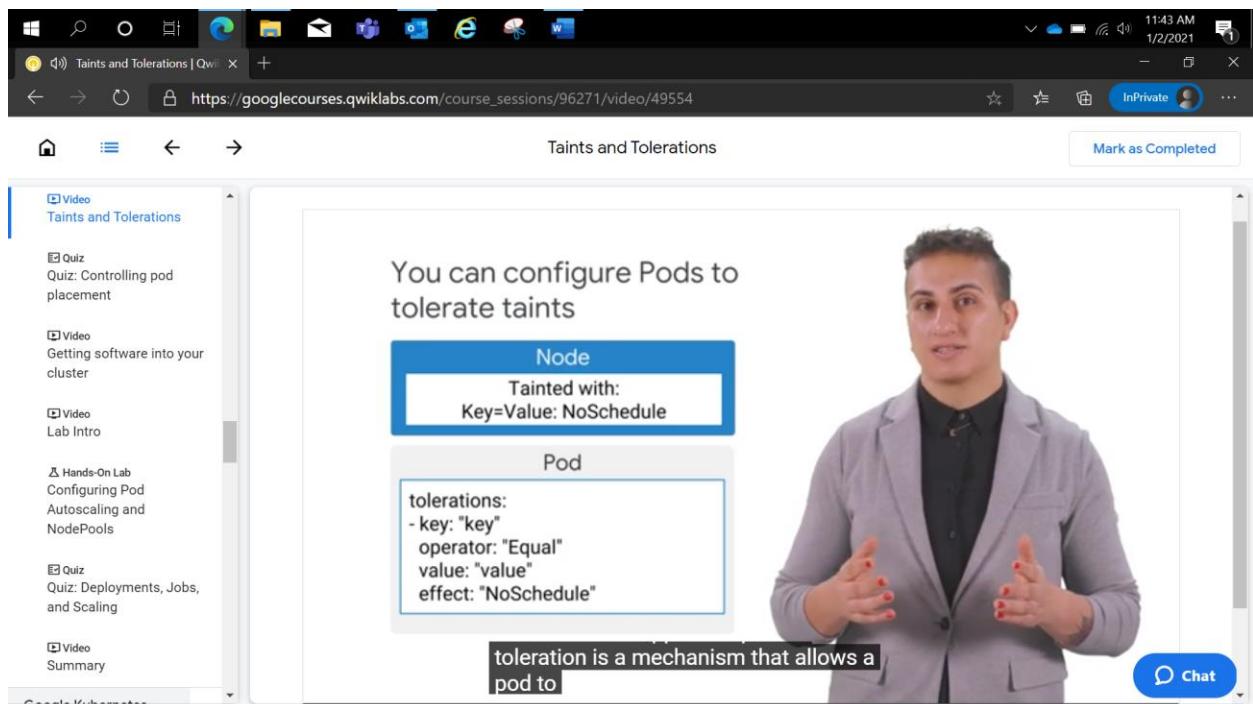
You can configure Pods to tolerate taints

Node
Tainted with:
Key=Value: NoSchedule

Pod
tolerations:
- key: "key"
operator: "Equal"
value: "value"
effect: "NoSchedule"

toleration is a mechanism that allows a pod to

Chat



11:44 AM 1/2/2021

Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Mark as Completed

InPrivate

Home Back Forward Taints and Tolerations

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Google Kubernetes

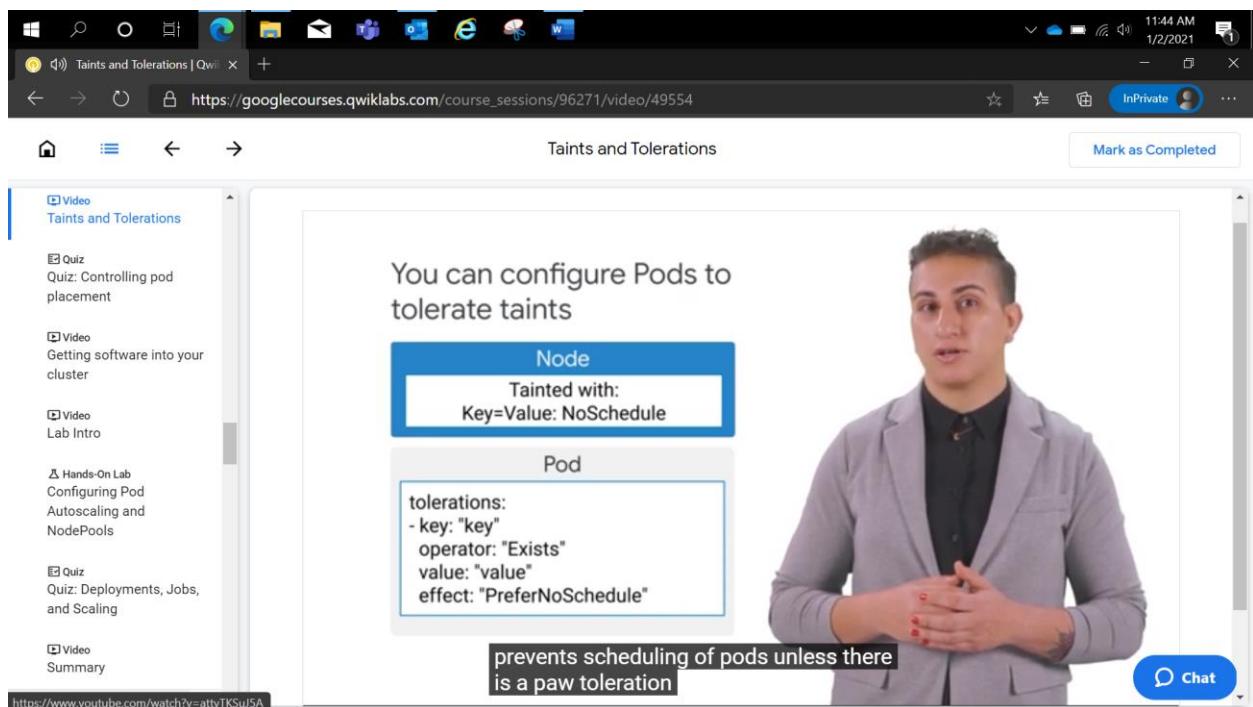
You can configure Pods to tolerate taints

Node
Tainted with:
Key=Value: NoSchedule

Pod
tolerations:
- key: "key"
operator: "Exists"
value: "value"
effect: "PreferNoSchedule"

prevents scheduling of pods unless there is a paw toleration

Chat



Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Taints and Tolerations

Mark as Completed

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Use node pools to manage different kinds of nodes

Cluster

Node pool "High CPU"

n1-highcpu-8 n1-highcpu-8
n1-highcpu-8 n1-highcpu-8

Node pool "High Memory"

n1-highmem-16 n1-highmem-16 n1-highmem-16

virtual machines have 16 virtual cpus and 104

Chat

Taints and Tolerations | Qwiklabs

https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554

Taints and Tolerations

Mark as Completed

Video Taints and Tolerations

Quiz Quiz: Controlling pod placement

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Direct Pods to desired nodes using nodeSelectors

```
apiVersion: v1
kind: Pod
metadata:
  name: nginx
spec:
  containers:
    - name: nginx
      image: nginx
  nodeSelector:
    cloud.google.com/gke-nodepool: small
[...]
```

each node pool with the node pool name you supply you

Chat

The screenshot shows a Microsoft Edge browser window with the URL https://googlecourses.qwiklabs.com/course_sessions/96271/video/49554. The title bar says "Taints and Tolerations | Qwiklabs". The main content area displays a video titled "Taints and Tolerations". A sidebar on the left lists course content: "Quiz: Controlling pod placement", "Video: Getting software into your cluster", "Video: Lab Intro", "Hands-On Lab: Configuring Pod Autoscaling and NodePools", "Quiz: Deployments, Jobs, and Scaling", and "Video: Summary". A code snippet in the center shows a YAML configuration for a Pod:

```
apiVersion: v1
kind: Pod
metadata:
  name: rendering-engine
spec:
  affinity:
    nodeAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        nodeSelectorTerms:
          - matchExpressions:
              - key: cloud.google.com/gke-nodepool
                operator: NotIn
                values:
                  - small
```

A callout box below the code states: "if node selectors are not expressive enough you can also use". A "Chat" button is in the bottom right.

The screenshot shows a Microsoft Edge browser window with the URL https://googlecourses.qwiklabs.com/course_sessions/96271/quizzes/49555. The title bar says "Quiz: Controlling pod placement | Qwiklabs". The main content area displays a quiz titled "Quiz: Controlling pod placement". A sidebar on the left lists course content: "Quiz: Controlling pod placement", "Video: Getting software into your cluster", "Video: Lab Intro", "Hands-On Lab: Configuring Pod Autoscaling and NodePools", "Quiz: Deployments, Jobs, and Scaling", and "Video: Summary". The quiz has two questions:

- True or false: if the labels on nodes change, affinity and anti-affinity rules will be applied to already-running Pods.
 - False
 - TrueA callout box says "That is correct."
- What is the difference between the taint system and the affinity system?
 - Taints apply to Pods, while affinity rules apply to nodes.
 - Taints are more modern and flexible than affinity rules, and should be used instead.
 - Taints apply to nodes, while affinity rules apply to Pods.A callout box says "That is correct."

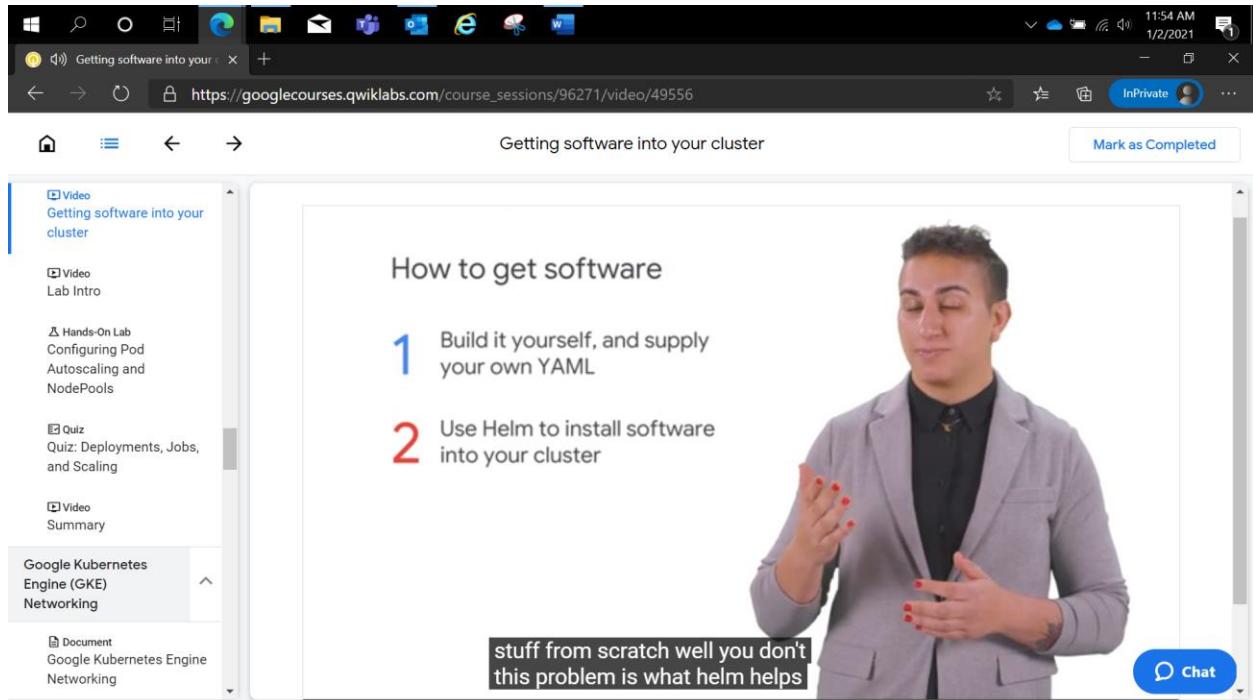
A "Chat" button is in the bottom right.

Getting software into your cluster

How to get software

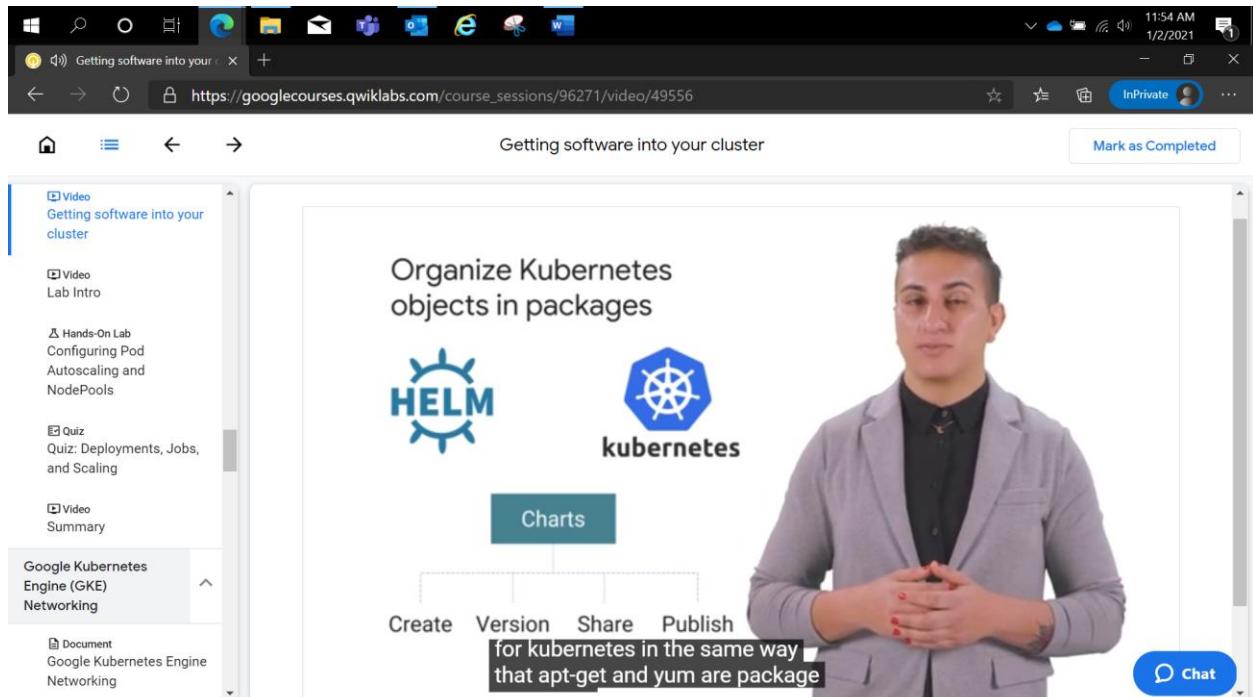
- 1 Build it yourself, and supply your own YAML
- 2 Use Helm to install software into your cluster

stuff from scratch well you don't this problem is what helm helps



Getting software into your cluster

Organize Kubernetes objects in packages



```
graph TD; Charts[Charts] --> Create[Create]; Charts --> Version[Version]; Charts --> Share[Share]; Charts --> Publish[Publish]
```

for kubernetes in the same way that apt-get and yum are package

Getting software into your cluster

Mark as Completed

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Google Kubernetes Engine (GKE) Networking

Document Google Kubernetes Engine Networking

Deploy complex packages

HELM

kubernetes

Charts

Download Configure Release

you need what resource constraints should they

Chat

Getting software into your cluster

Mark as Completed

Video Getting software into your cluster

Video Lab Intro

Hands-On Lab Configuring Pod Autoscaling and NodePools

Quiz Quiz: Deployments, Jobs, and Scaling

Video Summary

Google Kubernetes Engine (GKE) Networking

Document Google Kubernetes Engine Networking

Helm architecture

Helm client (helm)

Helm server (Tiller)

Kubernetes APIserver

Kubernetes cluster

it also stores the objects that represent helm chart releases

Chat

Getting software into your cluster

How to get software

- 1 Build it yourself, and supply your own YAML.
- 2 Use Helm to install software into your cluster.
- 3 Use Google Cloud Marketplace to install both open-source and commercial software.

even simpler
gcp marketplace

Configuring Pod Autoscaling and NodePools

1 hourFree

Rate Lab

Overview

In this lab, you set up an application in Google Kubernetes Engine (GKE), and then use a HorizontalPodAutoscaler to autoscale the web application. You then work with multiple node pools of different types, and you apply taints and tolerations to control the scheduling of Pods with respect to the underlying node pool.

Objectives

In this lab, you learn how to perform the following tasks:

- Configure autoscaling and HorizontalPodAutoscaler
- Add a node pool and configure taints on the nodes for Pod anti-affinity
- Configure an exception for the node taint by adding a toleration to a Pod's manifest

Task 0. Lab Setup

Access Qwiklabs

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an **incognito window**.

02:00:00

2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

START LAB

3. When ready, click **START LAB**.
4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

Open Google Console

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more](#).

Username

google2876526_student@qwiklabs.n



Password

TG959yrKDX



GCP Project ID

qwiklabs-gcp-0855e773352d3560



[New to labs? View our introductory video!](#)

5. Click **Open Google Console**.
6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or **incur charges**.

7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it.

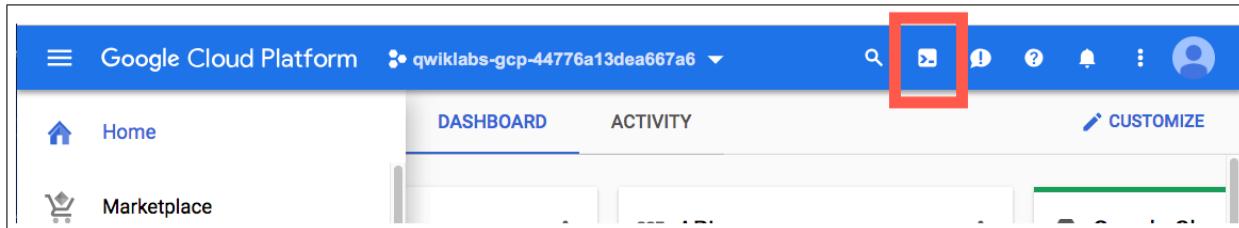
This clears your work and removes the project.

After you complete the initial sign-in steps, the project dashboard appears.

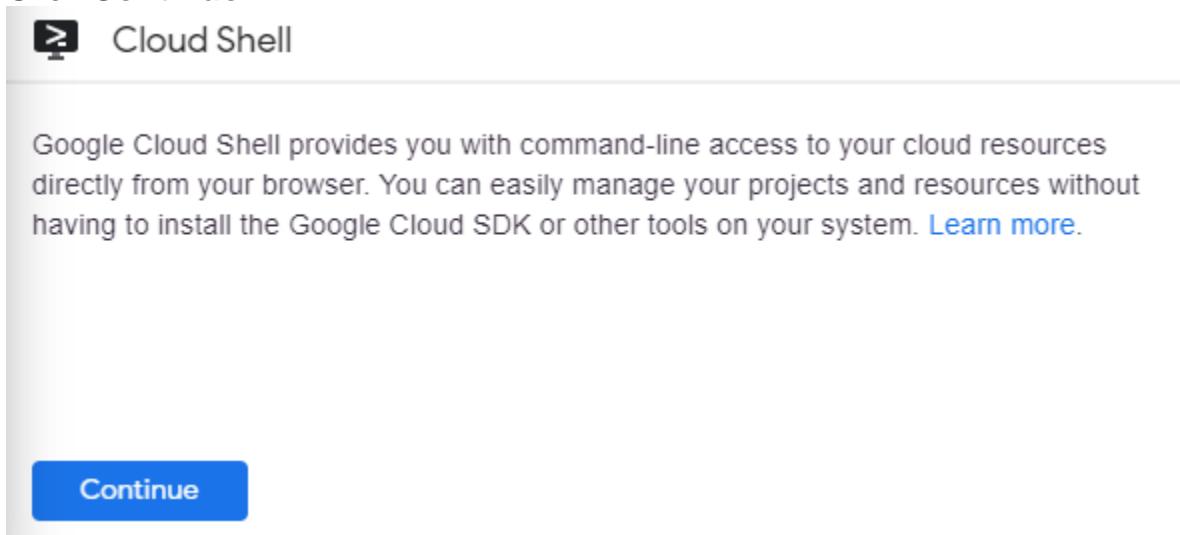
Activate Google Cloud Shell

Google Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Google Cloud Shell provides command-line access to your GCP resources.

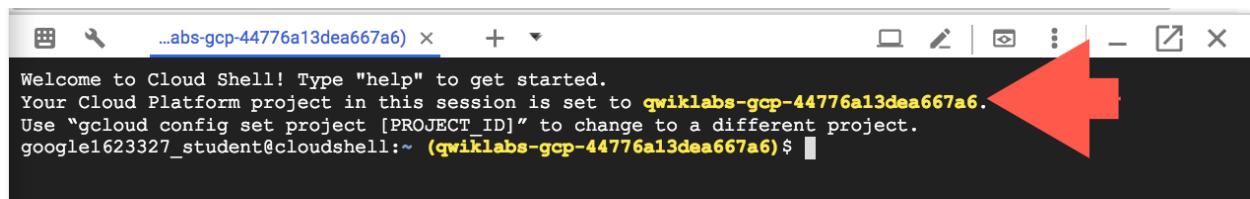
1. In GCP console, on the top right toolbar, click the Open Cloud Shell button.



2. Click **Continue**.



It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your *PROJECT_ID*. For example:



A screenshot of the Cloud Shell terminal window. The title bar shows "...abs-gcp-44776a13dea667a6". The terminal output shows: "Welcome to Cloud Shell! Type "help" to get started." "Your Cloud Platform project in this session is set to **qwiklabs-gcp-44776a13dea667a6**." "Use "gcloud config set project [PROJECT_ID]" to change to a different project." "google1623327_student@cloudshell:~ (qwiklabs-gcp-44776a13dea667a6)\$". A red arrow points to the project ID in the output.

gcloud is the command-line tool for Google Cloud Platform. It comes pre-installed on Cloud Shell and supports tab-completion.

You can list the active account name with this command:

```
gcloud auth list  
content_copy
```

Output:

```
Credentialed accounts:  
- <myaccount>@<mydomain>.com (active) content_copy
```

Example output:

```
Credentialed accounts:  
- google1623327 student@qwiklabs.netcontent_copy
```

You can list the project ID with this command:

```
gcloud config list project
```

```
content_copy
```

Output:

```
[core]
project = <project_ID>content_copy
```

Example output:

```
[core]
project = qwiklabs-gcp-44776a13dea667a6content_copy
```

Full documentation of **gcloud** is available on [Google Cloud gcloud Overview](#).

Task 1. Connect to the lab GKE cluster and deploy a sample workload

In this task, you connect to the lab GKE cluster and create a deployment manifest for a set of Pods within the cluster.

Connect to the lab GKE cluster

1. In Cloud Shell, type the following command to set the environment variable for the zone and cluster name.

```
export my_zone=us-central1-a
export my_cluster=standard-cluster-1
content_copy
```

2. Configure tab completion for the kubectl command-line tool.

```
source <(kubectl completion bash)
content_copy
```

3. Configure access to your cluster for kubectl:

```
gcloud container clusters get-credentials $my_cluster --zone $my_zone
content_copy
```

Deploy a sample web application to your GKE cluster

You will deploy a sample application to your cluster using the `web.yaml` deployment file that has been created for you:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: web
spec:
  replicas: 1
  selector:
    matchLabels:
      run: web
  template:
    metadata:
      labels:
        run: web
    spec:
      containers:
        - image: gcr.io/google-samples/hello-app:1.0
          name: web
          ports:
            - containerPort: 8080
              protocol: TCP
content_copy
```

This manifest creates a deployment using a sample web application container image that listens on an HTTP server on port 8080.

1. In Cloud Shell enter the following command to clone the repository to the lab Cloud Shell.

```
git clone https://github.com/GoogleCloudPlatform/training-data-analyst
content_copy
```

2. Create a soft link as a shortcut to the working directory.

```
ln -s ~/training-data-analyst/courses/ak8s/v1.1 ~/ak8s
content_copy
```

3. Change to the directory that contains the sample files for this lab.

```
cd ~/ak8s/Autoscaling/
content_copy
```

4. To create a deployment from this file, execute the following command:

```
kubectl create -f web.yaml --save-config
content_copy
```

5. Create a service resource of type NodePort on port 8080 for the web deployment.

```
kubectl expose deployment web --target-port=8080 --type=NodePort  
content_copy
```

6. Verify that the service was created and that a node port was allocated:

```
kubectl get service web  
content_copy
```

Your IP address and port number might be different from the example output.

Output (do not copy)

NAME	TYPE	CLUSTER-IP	EXTERNAL-IP	PORT(S)	AGE
web	NodePort	10.11.246.185	<none>	8080:32056/TCP	12s
content_copy					

Click *Check my progress* to verify the objective.

Deploy a sample web application to GKE cluster

Check my progress

Task 2. Configure autoscaling on the cluster

In this task, you configure the cluster to automatically scale the sample application that you deployed earlier.

Configure autoscaling

1. Get the list of deployments to determine whether your sample web application is still running.

```
kubectl get deployment  
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
web	1/1	1	1	5m48s
content_copy				

If the web deployment of your application is not displayed, return to task 1 and redeploy it to the cluster.

2. To configure your sample application for autoscaling (and to set the maximum number of replicas to four and the minimum to one, with a CPU utilization target of 1%), execute the following command:

```
kubectl autoscale deployment web --max 4 --min 1 --cpu-percent 1  
content_copy
```

When you use `kubectl autoscale`, you specify a maximum and minimum number of replicas for your application, as well as a CPU utilization target.

3. Get the list of deployments to verify that there is still only one deployment of the web application.

```
kubectl get deployment  
content_copy
```

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
web	1/1	1	1	8m21s
content_copy				

Inspect the HorizontalPodAutoscaler object

The `kubectl autoscale` command you used in the previous task creates a `HorizontalPodAutoscaler` object that targets a specified resource, called the `jscale` target, and scales it as needed. The autoscaler periodically adjusts the number of replicas of the scale target to match the average CPU utilization that you specify when creating the autoscaler.

1. To get the list of HorizontalPodAutoscaler resources, execute the following command:

```
kubectl get hpa  
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
web	Deployment/web	0%/1%	1	4	1	1m
content_copy						

2. To inspect the configuration of HorizontalPodAutoscaler in YAML form, execute the following command:

```
kubectl describe horizontalpodautoscaler web  
content_copy
```

The output should look like the example.

Output (do not copy)

```
Name: web  
Namespace: default  
Labels: <none>  
Annotations: <none>  
CreationTimestamp: Tue, 08 Sep 2020...  
Reference: Deployment/web  
Metrics:  
  resource cpu on pods (as a percentage of request): 0% (0) / 1%  
Min replicas: 1  
Max replicas: 4  
Deployment pods: 1 current / 1 desired  
Conditions:  
  Type Status Reason Message  
  ---- ---- - ----  
  AbleToScale True ScaleDownStabilized recent recommendations [...]  
  ScalingActive True ValidMetricFound the HPA was able to [...]  
  ScalingLimited False DesiredWithinRange the desired count [...]  
Events: <none>  
content_copy
```

3. To view the configuration of HorizontalPodAutoscaler in YAML form, execute the following command:

```
kubectl get horizontalpodautoscaler web -o yaml  
content_copy
```

The output should look like the example.

Output (do not copy)

```
apiVersion: autoscaling/v1
kind: HorizontalPodAutoscaler
metadata:
  annotations:
    autoscaling.alpha.kubernetes.io/conditions: [...]
    autoscaling.alpha.kubernetes.io/current-metrics: [...]
  creationTimestamp: 2018-11-14T02:59:28Z
  name: web
  namespace: default
  resourceVersion: "14588"
  selfLink: /apis/autoscaling/v1/namespaces/[...]
spec:
  maxReplicas: 4
  minReplicas: 1
  scaleTargetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: web
  targetCPUUtilizationPercentage: 1
status:
  currentCPUUtilizationPercentage: 0
  currentReplicas: 1
  desiredReplicas: 1
content_copy
```

Test the autoscale configuration

You need to create a heavy load on the web application to force it to scale out. You create a configuration file that defines a deployment of four containers that run an infinite loop of HTTP queries against the sample application web server.

You create the load on your web application by deploying the loadgen application using the `loadgen.yaml` file that has been provided for you.

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: loadgen
spec:
  replicas: 4
  selector:
    matchLabels:
      app: loadgen
  template:
    metadata:
      labels:
        app: loadgen
    spec:
      containers:
```

```
- name: loadgen
  image: k8s.gcr.io/busybox
  args:
    - /bin/sh
    - -c
    - while true; do wget -q -O- http://web:8080; done
content_copy
```

1. To deploy this container, execute the following command:

```
kubectl apply -f loadgen.yaml
content_copy
```

After you deploy this manifest, the web Pod should begin to scale.

Click *Check my progress* to verify the objective.

Deploying the loadgen application

Check my progress

2. Get the list of deployments to verify that the load generator is running.

```
kubectl get deployment
content_copy
```

The output should look like the example.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
loadgen	4/4	4	4	26s
web	1/1	1	1	14m

```
content_copy
```

3. Inspect HorizontalPodAutoscaler.

```
kubectl get hpa
content_copy
```

Once the loadgen Pod starts to generate traffic, the web deployment CPU utilization begins to increase. In the example output, the targets are now at 35% CPU utilization compared to the 1% CPU threshold.

Output (do not copy)

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
web	Deployment/web	35%/1%	1	4	1	8m

```
content_copy
```

4. After a few minutes, inspect the HorizontalPodAutoscaler again.

```
kubectl get hpa
content_copy
```

The autoscaler has increased the web deployment to four replicas.

Output (do not copy)

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
web	Deployment/web	88%/1%	1	4	4	9m
content_copy						

5. To stop the load on the web application, scale the loadgen deployment to zero replicas.

```
kubectl scale deployment loadgen --replicas 0  
content_copy
```

6. Get the list of deployments to verify that loadgen has scaled down.

```
kubectl get deployment  
content_copy
```

The loadgen deployment should have zero replicas.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
loadgen	0/0	0	0	3m40s
web	2/4	4	2	18m
web	2/4	4	2	18m
content_copy				

Note: You need to wait 2 to 3 minutes before you list the deployments again.

7. Get the list of deployments to verify that the web application has scaled down to the minimum value of 1 replica that you configured when you deployed the autoscaler.

```
kubectl get deployment  
content_copy
```

You should now have one deployment of the web application.

Output (do not copy)

NAME	READY	UP-TO-DATE	AVAILABLE	AGE
loadgen	0/0	0	0	6m28s
web	4/4	4	4	20m
content_copy				

Task 3. Manage node pools

In this task, you create a new pool of nodes using preemptible instances, and then you constrain the web deployment to run only on the preemptible nodes.

Add a node pool

1. To deploy a new node pool with three preemptible VM instances, execute the following command:

```
gcloud container node-pools create "temp-pool-1" \
--cluster=$my_cluster --zone=$my_zone \
--num-nodes "2" --node-labels=temp=true --preemptible
content_copy
```

If you receive an error that no preemptible instances are available you can remove the `--preemptible` option to proceed with the lab.

2. Get the list of nodes to verify that the new nodes are ready.

```
kubectl get nodes
content_copy
```

You should now have 4 nodes.

Your names will be different from the example output.

Output (do not copy)

NAME	STATUS	ROLES	AGE	VERSION
gke-standard-cluster-1-default-pool...xc	Ready	<none>	33m	v1.15.12-
gke.2				
gke-standard-cluster-1-default-pool...q8	Ready	<none>	33m	v1.15.12-
gke.2				
gke-standard-cluster-1-temp-pool-1...vj	Ready	<none>	32s	v1.15.12-
gke.2				
gke-standard-cluster-1-temp-pool-1...xj	Ready	<none>	37s	v1.15.12-
gke.2				
content_copy				

All the nodes that you added have the `temp=true` label because you set that label when you created the node-pool. This label makes it easier to locate and configure these nodes.

- To list only the nodes with the `temp=true` label, execute the following command:

```
kubectl get nodes -l temp=true  
content_copy
```

You should see only the two nodes that you added.

Your names will be different from the example output.

Output (do not copy)

NAME	STATUS	ROLES	AGE	VERSION
gke-standard-cluster-1-temp-pool-1-...vj v1.15.12-gke.2	Ready	<none>	3m26s	
gke-standard-cluster-1-temp-pool-1-...xj v1.15.12-gke.2	Ready	<none>	3m31s	

Control scheduling with taints and tolerations

To prevent the scheduler from running a Pod on the temporary nodes, you add a taint to each of the nodes in the temp pool. Taints are implemented as a key-value pair with an effect (such as `NoExecute`) that determines whether Pods can run on a certain node. Only nodes that are configured to tolerate the key-value of the taint are scheduled to run on these nodes.

- To add a taint to each of the newly created nodes, execute the following command.

You can use the `temp=true` label to apply this change across all the new nodes simultaneously.

```
kubectl taint node -l temp=true nodetype=preemptible:NoExecute  
content_copy
```

To allow application Pods to execute on these tainted nodes, you must add a `tolerations` key to the deployment configuration.

- Edit the `web.yaml` file to add the following key in the template's `spec` section:

```
tolerations:  
- key: "nodetype"
```

```
operator: Equal
value: "preemptible"
content_copy
```

The spec section of the file should look like the example.

```
...
spec:
  tolerations:
  - key: "nodetype"
    operator: Equal
    value: "preemptible"
  containers:
  - image: gcr.io/google-samples/hello-app:1.0
    name: web
    ports:
    - containerPort: 8080
      protocol: TCP
content_copy
```

3. To force the web deployment to use the new node-pool add a `nodeSelector` key in the template's `spec` section. This is parallel to the `tolerations` key you just added.

```
nodeSelector:
  temp: "true"
content_copy
```

Note: GKE adds a custom label to each node called `cloud.google.com/gke-nodepool` that contains the name of the node-pool that the node belongs to. This key can also be used as part of a `nodeSelector` to ensure Pods are only deployed to suitable nodes.

The full `web.yaml` deployment should now look as follows.

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: web
spec:
  replicas: 1
  selector:
    matchLabels:
      run: web
  template:
    metadata:
      labels:
        run: web
    spec:
      tolerations:
      - key: "nodetype"
        operator: Equal
        value: "preemptible"
      nodeSelector:
        temp: "true"
      containers:
```

```
- image: gcr.io/google-samples/hello-app:1.0
  name: web
  ports:
  - containerPort: 8080
    protocol: TCP
content_copy
```

4. To apply this change, execute the following command:

```
kubectl apply -f web.yaml
content_copy
```

If you have problems editing this file successfully you can use the pre-prepared sample file called `web-tolerations.yaml` instead.

Click *Check my progress* to verify the objective.

Manage node pools

Check my progress

5. Get the list of Pods.

```
kubectl get pods
content_copy
```

Your names might be different from the example output.

Output (do not copy)

NAME	READY	STATUS	RESTARTS	AGE
web-7cb566bccd-pkfst	1/1	Running	0	1m

```
content_copy
```

6. To confirm the change, inspect the running web Pod(s) using the following command

```
kubectl describe pods -l run=web
content_copy
```

A Tolerations section with `nodetype=preemptible` in the list should appear near the bottom of the (truncated) output.

Output (do not copy)

```
<SNIP>

Node-Selectors:  temp=true
Tolerations:    node.kubernetes.io/not-ready:NoExecute op=Exists for 300s
                node.kubernetes.io/unreachable:NoExecute op=Exists for 300s
                nodetype=preemptible
Events:
```

```
<SNIP>
content_copy
```

The output confirms that the Pods will tolerate the taint value on the new preemptible nodes, and thus that they can be scheduled to execute on those nodes.

7. To force the web application to scale out again scale the loadgen deployment back to four replicas.

```
kubectl scale deployment loadgen --replicas 4
content_copy
```

You could scale just the web application directly but using the loadgen app will allow you to see how the different taint, toleration and nodeSelector settings that apply to the web and loadgen applications affect which nodes they are scheduled on.

8. Get the list of Pods using the wide output format to show the nodes running the Pods

```
kubectl get pods -o wide
content_copy
```

This shows that the loadgen app is running only on `default-pool` nodes while the web app is running only the preemptible nodes in `temp-pool-1`.

The taint setting prevents Pods from running on the preemptible nodes so the loadgen application only runs on the default pool. The toleration setting allows the web application to run on the preemptible nodes and the nodeSelector forces the web application Pods to run on those nodes.

NAME	READY	STATUS	[...]	NODE
Loadgen-x0	1/1	Running	[...]	gke-xx-default-pool-y0
loadgen-x1	1/1	Running	[...]	gke-xx-default-pool-y2
loadgen-x3	1/1	Running	[...]	gke-xx-default-pool-y3
loadgen-x4	1/1	Running	[...]	gke-xx-default-pool-y4
web-x1	1/1	Running	[...]	gke-xx-temp-pool-1-z1
web-x2	1/1	Running	[...]	gke-xx-temp-pool-1-z2
web-x3	1/1	Running	[...]	gke-xx-temp-pool-1-z3
web-x4	1/1	Running	[...]	gke-xx-temp-pool-1-z4

Quiz: Deployments, Jobs, and Scaling

Your score: 58% Passing score: 91% [Retake](#)

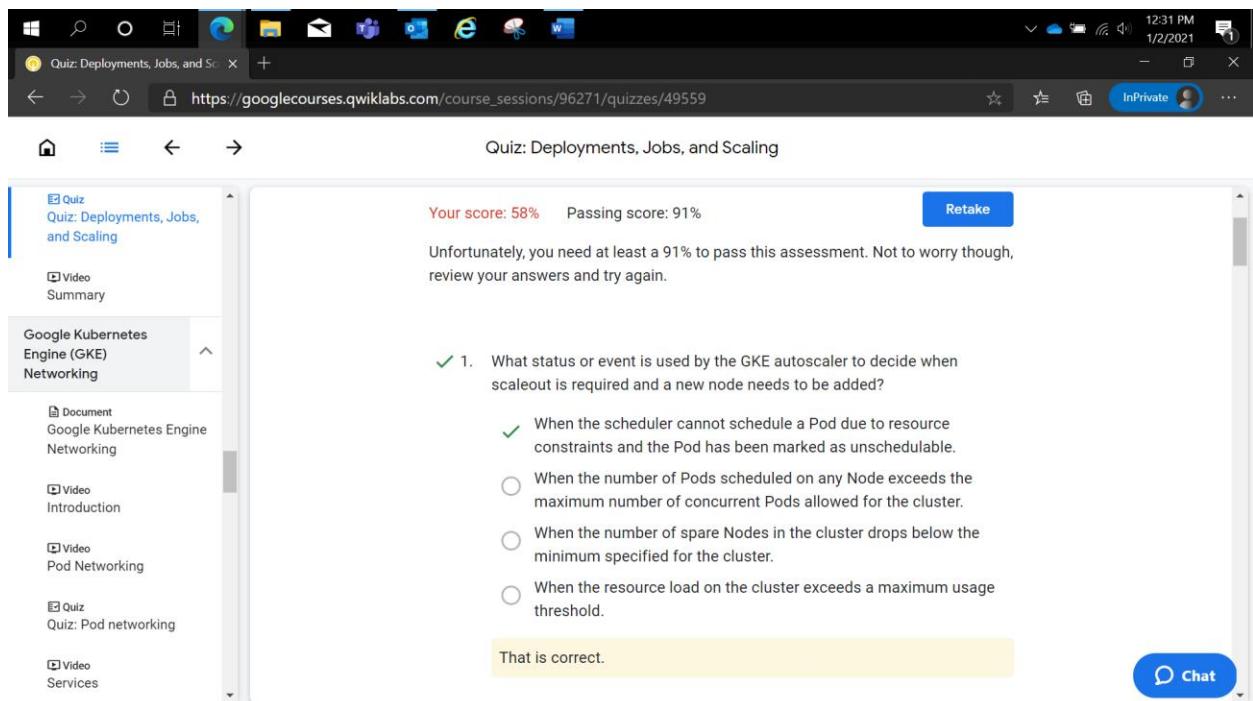
Unfortunately, you need at least a 91% to pass this assessment. Not to worry though, review your answers and try again.

✓ 1. What status or event is used by the GKE autoscaler to decide when scaleout is required and a new node needs to be added?

- When the scheduler cannot schedule a Pod due to resource constraints and the Pod has been marked as unschedulable.
- When the number of Pods scheduled on any Node exceeds the maximum number of concurrent Pods allowed for the cluster.
- When the number of spare Nodes in the cluster drops below the minimum specified for the cluster.
- When the resource load on the cluster exceeds a maximum usage threshold.

That is correct.

Chat



Quiz: Deployments, Jobs, and Scaling

✗ label: failure-domain.beta.kubernetes.io/zone

zone: failure-domain.beta.kubernetes.io/zone

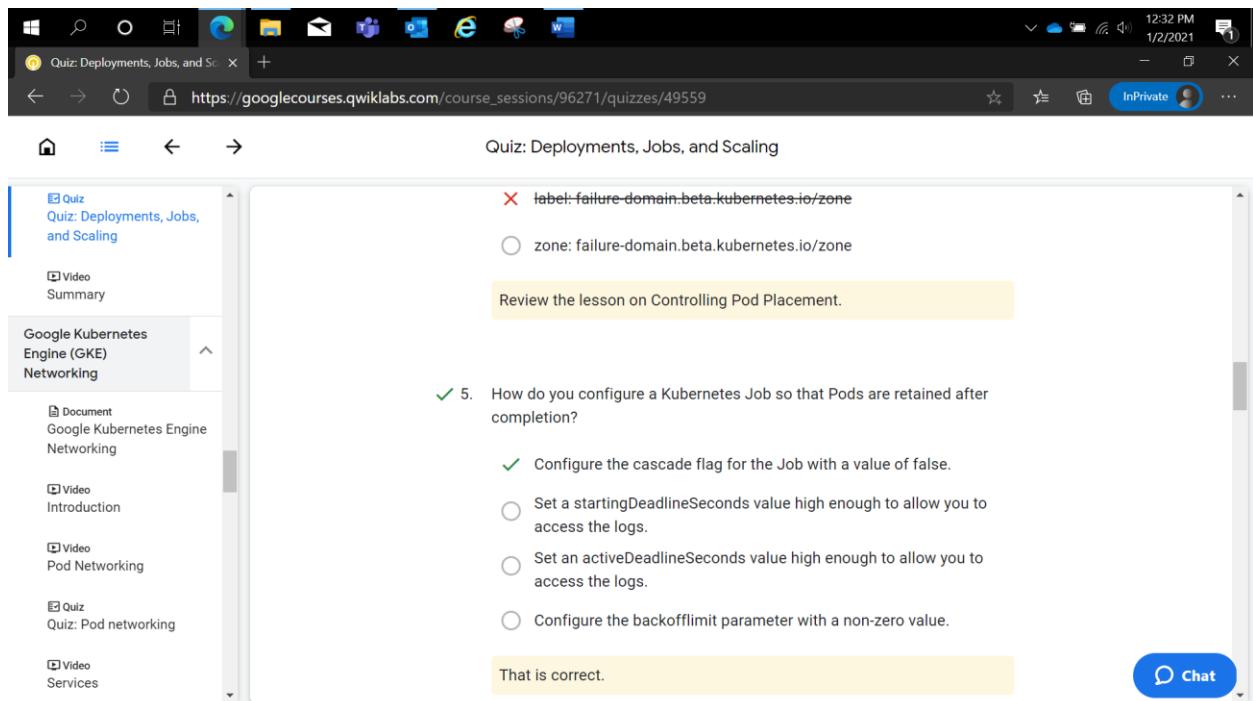
Review the lesson on Controlling Pod Placement.

✓ 5. How do you configure a Kubernetes Job so that Pods are retained after completion?

- Configure the cascade flag for the Job with a value of false.
- Set a startingDeadlineSeconds value high enough to allow you to access the logs.
- Set an activeDeadlineSeconds value high enough to allow you to access the logs.
- Configure the backofflimit parameter with a non-zero value.

That is correct.

Chat



Quiz: Deployments, Jobs, and Scaling

7. You are configuring the rollout strategy for your Deployment that contains 8 Pods. You need to specify a Deployment property that will ensure at least 75% of the desired number of Pods is always running at the same time. What property and value should you set for the deployment to ensure that this is the case?

maxUnavailable=25%

maxSurge=25%

maxSurge=2

maxUnavailable=2

That is correct.

8. You are resolving a range of issues with a Deployment and need to make a large number of changes. Which command can you execute to group these changes into a single rollout, thus avoiding pushing out a large

Quiz: Deployments, Jobs, and Scaling

8. You are resolving a range of issues with a Deployment and need to make a large number of changes. Which command can you execute to group these changes into a single rollout, thus avoiding pushing out a large number of rollouts?

kubectl rollout pause deployment

kubectl stop deployment

kubectl delete deployment

kubectl rollout resume deployment

That is correct.

9. You have configured a Kubernetes Job with a backofflimit of 4 and a completions count of 8. If the Pods launched by the Job continually fail,

Quiz: Deployments, Jobs, and Scaling

That is correct.

9. You have configured a Kubernetes Job with a backofflimit of 4 and a completions count of 8. If the Pods launched by the Job continually fail, how long does it take for four failures to happen and what does the Job report?

80 seconds, Job fails with BackoffLimitExceeded as the reason.
 40 seconds, Job continues launching Pods until completion reaches 8.
 80 seconds, Job continues launching Pods until completion reaches 8.
 40 seconds, Job fails with BackoffLimitExceeded as the reason.

That is correct.

✓ 10. You have made a number of changes to your deployment and applied those changes. Which command should you use to rollback the environment to the deployment identified in the deployment history as revision 2?

Quiz: Deployments, Jobs, and Scaling

That is correct.

10. You have made a number of changes to your deployment and applied those changes. Which command should you use to rollback the environment to the deployment identified in the deployment history as revision 2?

Run 'kubectl rollout undo deployment --to-revision=2'.
 Select the desired revision from the revision history list in the Google Cloud console.
 Run 'kubectl apply -f DEPLOYMENT_FILE --to-revision=2'.
 Run 'kubectl rollout undo deployment' twice.

That is correct.

Quiz: Deployments, Jobs, and Scaling

Run 'kubectl rollout undo deployment' twice.

That is correct.

11. After a Deployment has been created and its component Pods are running, which component is responsible for ensuring that a replacement Pod is launched whenever a Pod fails or is evicted?

ReplicaSet

DaemonSet

StatefulSet

Deployment

That is correct.

Chat

Quiz: Deployments, Jobs, and Scaling

Shut themselves down.

Review the Jobs and CronJobs lesson.

3. A parallel Kubernetes Job is configured with parallelism of property of 4 and a completion property of 9. How many Pods are kept in a running state by the Job controller immediately after the sixth successful completion?

6

1

3

4

That is correct.

Chat

Quiz: Deployments, Jobs, and Scaling

Set an activeDeadlineSeconds value high enough to allow you to access the logs.

That is correct.

6. You are configuring a Job to process the conversion of a sample of a large number of video files from one format to another. Which parameter should you configure to ensure that you stop processing once a sufficient quantity have been processed?

- replicas=4
- parallelism=4
- completions=4
- backofflimit=4

That is correct.

Chat

Quiz: Deployments, Jobs, and Scaling

That is correct.

7. You have autoscaling enabled on your cluster. What conditions are required for the autoscaler to decide to delete a node?

- If the overall cluster is underutilized, a randomly selected node is deleted.
- If a node is underutilized and there are no Pods currently running on the Node.
- If a node is underutilized and running Pods can be run on other Nodes.
- If the overall cluster is underutilized, the least busy node is deleted.

That is correct.

8. What status or event is used by the GKE autoscaler to decide when

Chat

Quiz: Deployments, Jobs, and Scaling

access the logs.

That is correct.

11. When specifying Inter-pod affinity rules, you need to specify an affinity rule at the zone level, not at the individual Node level. Which additional parameter in the Pod manifest YAML must you set to apply this override?

- matchLabels: failure-domain.beta.kubernetes.io/zone
- zone: failure-domain.beta.kubernetes.io/zone
- topologyKey: failure-domain.beta.kubernetes.io/zone
- label: failure-domain.beta.kubernetes.io/zone

That is correct.

Chat

Quiz: Deployments, Jobs, and Scaling

Your score: 100% Passing score: 91%

Congratulations! You passed this assessment.

Retake

1. With a Kubernetes Job configured with a parallelism value of 3 and no completion count what happens to the status of the Job when one of the Pods successfully terminates?

- The Job is not considered complete until all Pods terminate successfully and shut themselves down.
- The entire Job is considered complete and the remaining Pods are shut down.
- The Job is considered complete, but the remaining Pods are left to shut themselves down.
- Pods in a parallel Job must be able to detect when other Pods have completed and should terminate automatically.

That is correct.

Chat