



---

## Assignment-2: Scaling Services

---

### Instructions:

1. This is an individual homework. Each student must submit his/her solution to the problem through Assignment-2 dropbox.
2. All submissions are to be via the Quercus Submission System; no other way of submission is accepted.
3. All submissions must be submitted through the dropbox by the due date or by the closure date, Clouse Date allows 3 days for late submissions.

### Purpose and Objective

This assignment addresses distributed systems fault-tolerance, synchronization, data consistency, and service performance issues.

### Keywords

Cluster, K8s, container composition, container deployment, serverless computing, authentication, authorization, Borg.

### Resources

1. The Artifact Registry Documentation: <https://cloud.google.com/artifact-registry/docs>
2. Google Cloud Managed Service for Prometheus.
3. Cloud Monitoring Overview: <https://cloud.google.com/monitoring/docs/monitoring-overview> [Links to an external site.](#)
4. Example: Collecting Apache Web Server Metrics: <https://cloud.google.com/monitoring/monitor-compute-engine-virtual-machine> [Links to an external site.](#)
  - a. **Generate a traffic to a web server:** `timeout 120 bash -c -- 'while true; do curl localhost; sleep $((RANDOM % 4)) ; done'`
  - b. Bash
    - i. <https://ss64.com/bash/syntax.html>
    - ii. <https://ss64.com/bash/>
  - c. curl command, Linux Manual: <https://ss64.com/bash/curl.html>
  - d. timeout command, Linux Manual: <https://ss64.com/bash/timeout.html>
  - e. sleep command, Linux manual: <https://ss64.com/bash/sleep.html>
5. Google Compute Engine Metrics Collection (Ops Agent and Open Telemetry Protocol): <https://cloud.google.com/monitoring/agent/ops-agent> [Links to an external site.](#)
6. Platforms  
Monitoring: <https://cloud.google.com/stackdriver/docs/solutions/agents/ops-agent/third-party> [Links to an external site.](#)



1. Google implemented the Borg system to manage clusters of servers/cells, identify the basic three functionalities that Borg accomplishes. Explain how separation of duties/isolation is accomplished in Borg. Identify the main components of Borg and describe how these components define its architecture.

**[20 marks]**

2. Kubernetes is a portable, extensible, open source platform for managing containerized workloads and services, that facilitates both declarative configuration and automation. It has a large, rapidly growing ecosystem. Kubernetes services, support, and tools are widely available.
  - The name Kubernetes originates from Greek, meaning helmsman or pilot. K8s as an abbreviation results from counting the eight letters between the "K" and the "s". Google open-sourced the Kubernetes project in 2014. Kubernetes combines [over 15 years of Google's experience](#) running production workloads at scale with best-of-breed ideas and practices from the community.

2.1 Identify when would you deploy Kubernetes.

2.2 Identify the main components of Kubernetes.

**[20 marks]**




3. Using the Google Cloud platform, create a serverless web service that scales automatically. Observe the performance of the service at high demand to verify the autoscaling feature. [30 marks]

- a. Start a new project on Google Cloud.
- b. Name your project as q3assignment2

---


### New Project

---

 You have 20 projects remaining in your quota. Request an increase or delete projects. [Learn more](#)


[MANAGE QUOTAS](#)

Project name \*

My Project 

Project ID: midyear-nebula-406522. It cannot be changed later. [EDIT](#)

Location \*

 No organisation [BROWSE](#)

Parent organisation or folder

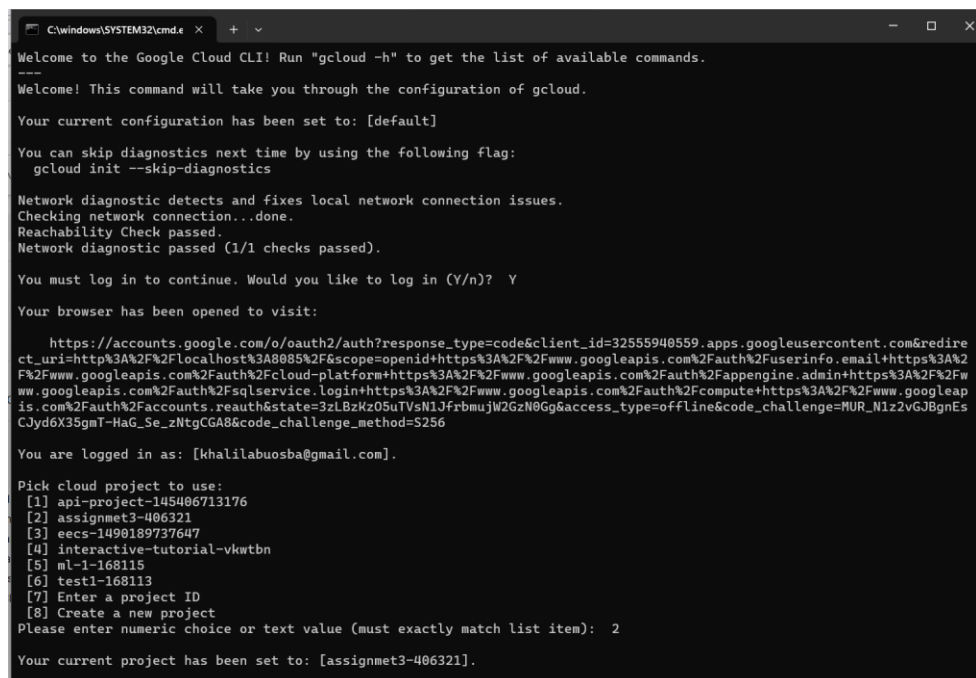
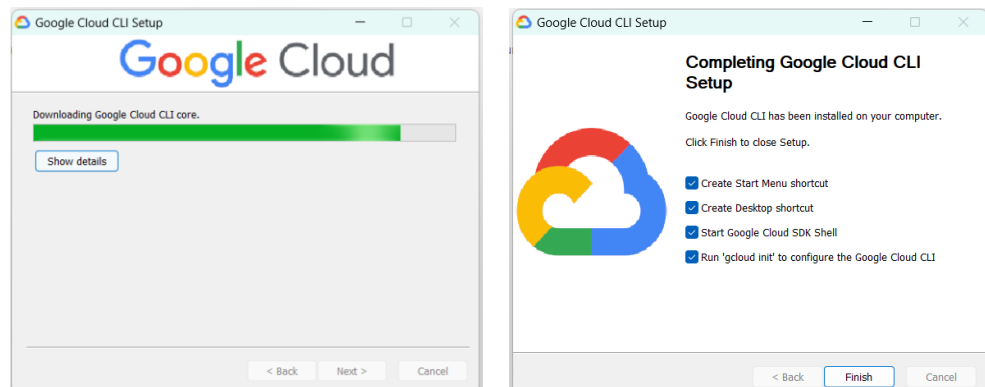
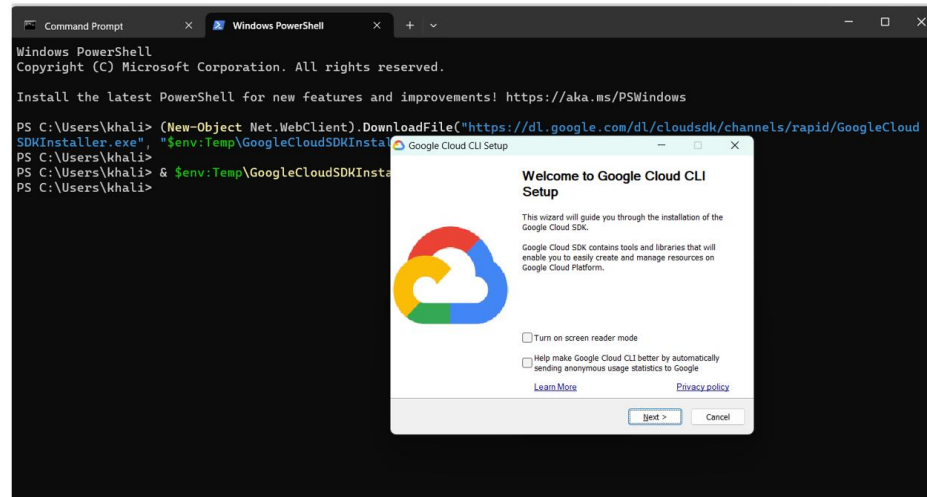
CREATE

CANCEL

- c. Install the Google Cloud Command Line Interface that matches your platform: <https://cloud.google.com/sdk/docs/install>
  - a. If your platform is Windows, start a PowerShell window and run the following command to download the CLI installer:  
**(New-Object Net.WebClient).DownloadFile("https://dl.google.com/dl/cloudsdk/channels/rapid/GoogleCloudSDKInstaller.exe", "\$env:Temp\GoogleCloudSDKInstaller.exe") & \$env:Temp\GoogleCloudSDKInstaller.exe**

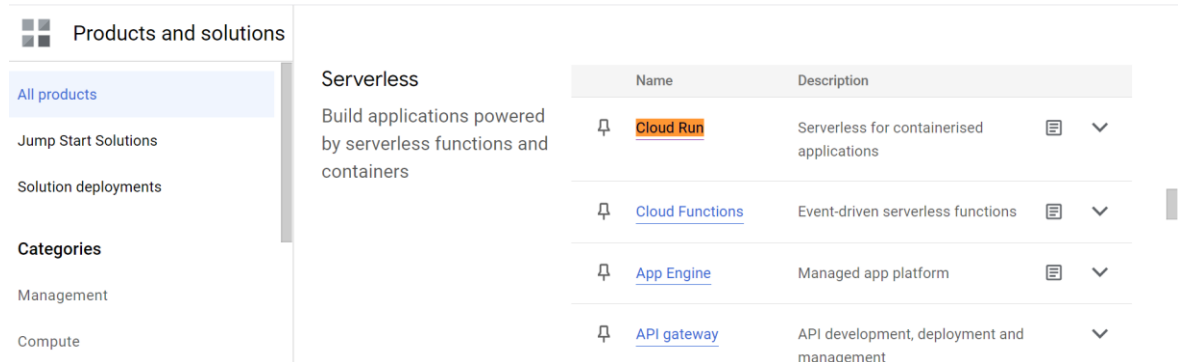


b.

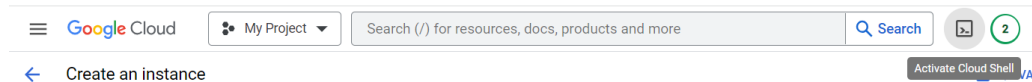




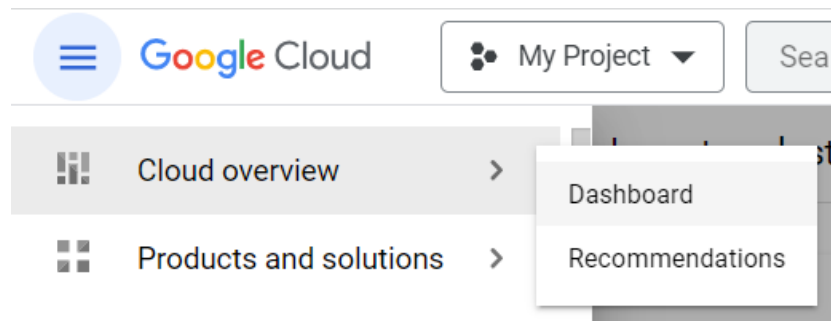
- d. Run the [CLI](#) initialization command and pick the project to use: **gcloud init**
- e. Under Serverless menu, of Products and Service, select Cloud Run:



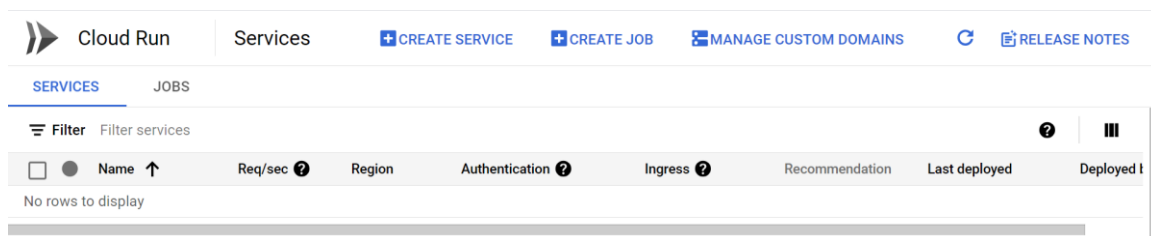
- f. Click > to activate the cloud shell



- g. To initialize the environment, list the projects and select your project, Run the command > **gcloud init**
- h. You can select the project using the command > **gcloud config set project\_name**
- i. Click on → Navigation Menu → Dashboard

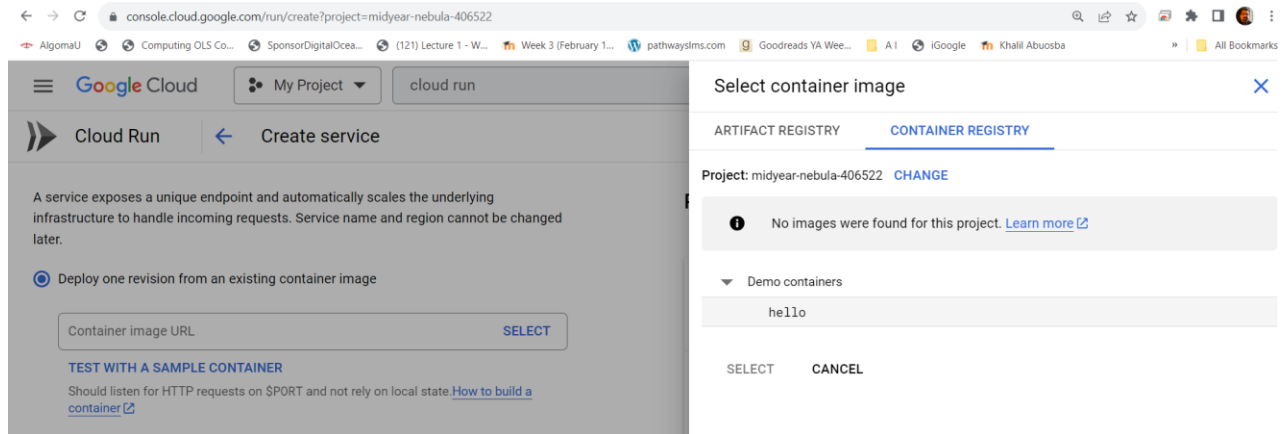


- j. Click → Cloud Run





- k. Click Select → Container Registry → Demo Containers → hello  
→ Select



- l. Configure the environment as the following except minimum number of instances as 1 instead of 3, click → Create

**CPU allocation and pricing**

☐ CPU is only allocated during request processing  
You are charged per request and only when the container instance processes a request.

☒ CPU is always allocated  
You are charged for the entire lifecycle of the container instance.

**Auto-scaling**

Minimum and maximum numbers of instances that the created revision scales to.

Minimum number of instances \* 3 Maximum number of instances \* 100

**Ingress control**

☐ Internal  
Allow traffic from your project, shared VPC and VPC service controls perimeter. Traffic from another Cloud Run service must be routed through a VPC. Limitations apply. [Learn more](#)

☒ All  
Allow direct access to your service from the Internet

**Authentication \***

☒ Allow unauthenticated invocations  
Tick this if you are creating a public API or website.

☐ Require authentication  
Manage authorised users with Cloud IAM.

**Container(s), volumes, networking, security**

**CREATE** CANCEL

**Pricing summary**

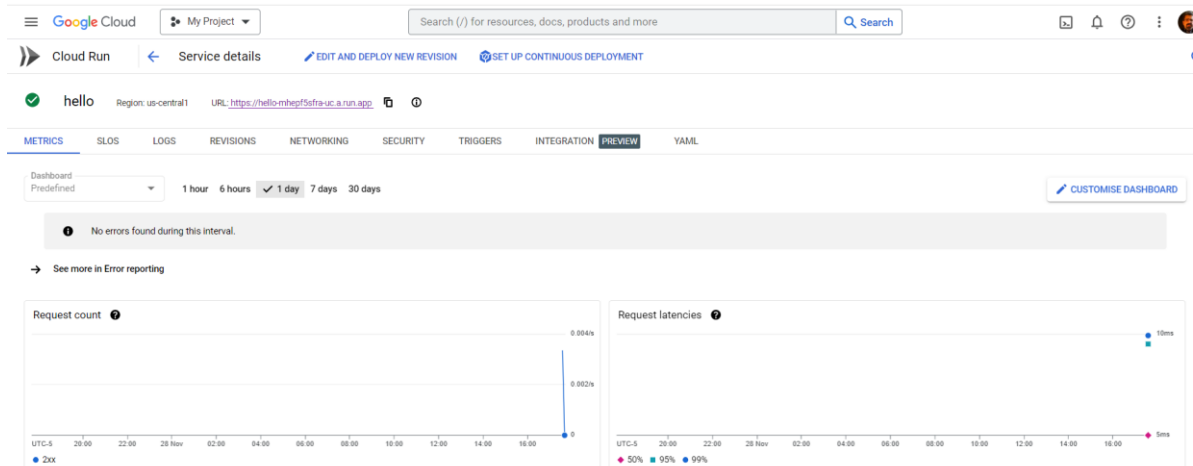
Cloud Run pricing

Free tier

First 240,000 vCPU-seconds/month  
First 450,000 GiB-seconds/month

→ Check paid tiers details

- m. From the YAML file, observe the number of replicas; how many replicas were specified?
- n. What is the URL specified for consuming the service.
- o. What is the result of parsing the URL in the browser.
- p. Click on View on Cloud Council:

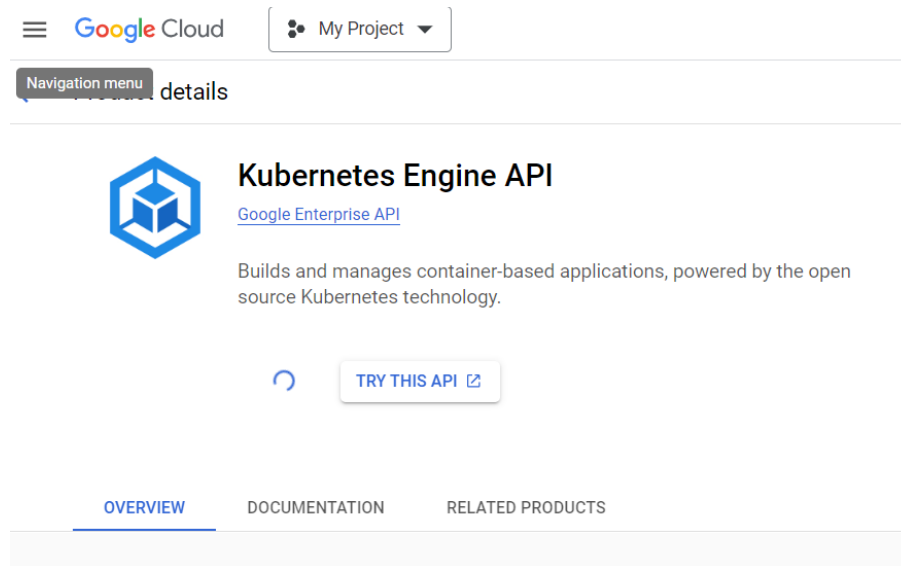


- q. Attempt to flood the service with http requests for 10 to 30 minutes (i.e. you may use curl, open multiple terminal or cmd sessions and request the contents of the homepage using the redundant curl commands).
- r. What is the analytical summary output of the Container Instance Count and container CPU Utilization?



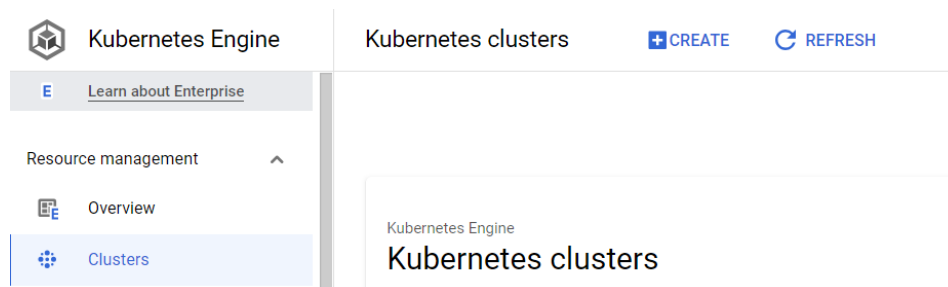
**4. Using the Google Kubernetes Engine (GKE) services, create a K8s cluster and deploy a service to the cluster. [30 marks]**

1. Make sure that the GKE API is enabled → Navigation Menu → Kubernetes Engine → Kubernetes API → Enable.



- a. Review the Create Cluster | Deploy Container tutorial:  
<https://console.cloud.google.com/kubernetes/list/overview?project=midyear-nebula-406522>
2. Create a Kubernetes *cluster*, which provides compute, storage, networking, and other services for applications, similar to a virtual data center.

- a. Click Kubernetes Engine → Cluster → Create







b. Specify the name of the cluster as autopilot-cluster-100 → Next

### Cluster basics

Create an Autopilot cluster by specifying a name and region. After the cluster has been created, you can deploy your workload through Kubernetes and we'll take care of the rest, including:

- ✓ **Nodes:** Automated node provisioning, scaling and maintenance
- ✓ **Networking:** VPC-native traffic routing for public or private clusters
- ✓ **Security:** Shielded GKE nodes and Workload Identity
- ✓ **Telemetry:** Cloud Operations logging and monitoring

Name

autopilot-cluster-100

Cluster names must start with a lowercase letter followed by up to 39 lowercase letters, numbers or hyphens. They can't end with a hyphen. You cannot change the cluster's name once it has been created.

Region

us-central1

The regional location in which your cluster's control plane and nodes are located. You cannot change the cluster's region once it has been created.

[NEXT: NETWORKING](#)

[RESET SETTINGS](#)

c. Next → Next → Next (Default Values) → Review and Create

d. → Create Cluster.

Kubernetes Engine

Kubernetes clusters

CREATE DEPLOY REFRESH ONBOARDING NEW OPERATIONS LEARN

OVERVIEW OBSERVABILITY COST OPTIMISATION

Filter Enter property name or value

Status	Name	Location	Mode	Number of nodes	Total vCPUs	Total memory	Notifications
<input checked="" type="checkbox"/>	autopilot-cluster-100	us-central1	Autopilot		0	0 GB	

e. Click observability.

b. Capture your version of the Container Error Logs/Second chart and paste in your report.

3. Service Composition:

a. This app (hello-app) composes a web server written in Go that responds to requests with the message Hello World!

b. The process requires the source code that may be imported from the github Google Cloud Platform repository using the command:



```
git clone https://github.com/GoogleCloudPlatform/kubernetes-engine-samples; cd
kubernetes-engine-samples/quickstarts/hello-app
```

- c. Set the `Project_Name` environment variable to your Google Cloud project name in the export command below.  
The `Project_Name` variable associates the container image with your project's Artifact Registry:  
`export PROJECT_ID = Your_Project_Name`

- d. Verify the project name using the echo command: `echo ${PROJECT_ID}`

e. **Build the Docker image**

- a. Create a repository (create a repository to store your container image)

1. specify a region for your repository. Run the following command to see a list of possible regions, then select the one closest to you:  
`gcloud artifacts locations list`
2. Create a **REGION** variable to hold your selected region. The following command sets `REGION` to `us-west1`:  
`export REGION=us-west1`
3. Run the following command to create your repository.  
`gcloud artifacts repositories create hello-repo -- repository-format=docker --location=${REGION} -- description="Docker repository"`

GKE accepts Docker images as the application deployment format. Before deploying the `hello-app` application, you must package the source code as a Docker image.

To build a Docker image, you need source code and a Dockerfile. A Dockerfile contains instructions on how the image is built.

- b. Build and tag the Docker image for `hello-app`:  
`docker build -t ${REGION}-docker.pkg.dev/${PROJECT_ID}/hello-repo/hello-app:v1 .`

```
kstarts/hello-app (midyear-nebula-406522) $ docker build -t us-west1-docker.pkg.dev/midyear-nebula-406522/hello-repo/hello-app:v1 .
[+] Building 36.7s (13/13) FINISHED
=> [internal] load build definition from Dockerfile
=> => transferring dockerfile: 983B
=> [internal] load .dockerignore
=> => transferring context: 2B
=> [internal] load metadata for gcr.io/distroless/base-debian11:latest
=> [internal] load metadata for docker.io/library/golang:1.21.0
=> [builder 1/5] FROM docker.io/library/golang:1.21.0@sha256:b490aef0eae153648dd3c5d25be59a63f966b5f9e1311245c947de4506981aa
=> => resolve docker.io/library/golang:1.21.0@sha256:b490aef0eae153648dd3c5d25be59a63f966b5f9e1311245c947de4506981aa
=> => sha256:b490aef0eae153648dd3c5d25be59a63f966b5f9e1311245c947de4506981aa 2.36kB / 2.36kB
=> => sha256:434a1905832963ddda9f85a8f329a0c0b78a3529c99d57259f0888e0feb973 7.06kB / 7.06kB
=> => sha256:de4cac68b6165c40cf6f8b30417948c31be03a968e233e55ee40221553a5e570 49.56MB / 49.56MB
```

This command instructs Docker to build the image using the `Dockerfile` in the current directory and tag it with a name, such as `us-west1-docker.pkg.dev/your-project_name/hello-repo/hello-app:v1`. The image is pushed to Artifact Registry in the next section

1. Run the `docker images` command to verify that the build was successful:  
`docker images`



## 2. Push the image to the repository

## 3. Enable the registry API:

**gcloud services enable artifactregistry.googleapis.com**

## 4. Configure the Docker command-line tool to authenticate to Artifact Registry:

**gcloud auth configure-docker \${REGION}-docker.pkg.dev**

## 5. Push the image to the repository

**docker push \${REGION}-docker.pkg.dev/\${PROJECT\_ID}/hello-repo/hello-app:v1**

```
khailaboosba@cloudshell:~/kubernetes-engine-samples/quickstarts/hello-app:~$ docker push $(REGION)-docker.pkg.dev/midyear-nebula-406522/hello-repo/hello-app:v1
The push refers to repository [us-west1-docker.pkg.dev/midyear-nebula-406522/hello-repo/hello-app]
9ac1c5cd32f7: Preparing
9ac1c5cd32f7: Pushed
714f56238fb5: Pushed
f33e343848bd: Pushed
4cb10dd2545b: Pushed
```

## 6. List all repositories

**gcloud artifacts repositories list**

```
kstarts/hello-app (midyear-nebula-406522)$ gcloud artifacts repositories list
Listing items under project midyear-nebula-406522, across all locations.

ARTIFACT_REGISTRY

REPOSITORY: hello-repo
FORMAT: DOCKER
MODE: STANDARD_REPOSITORY
DESCRIPTION: Docker repository
LOCATION: us-west1
LABELS:
ENCRYPTION: Google-managed key
CREATE_TIME: 2023-11-29T06:27:03
```

## 4. Deploy a sample app to the K8s cluster. Apps and their associated services that are running in Kubernetes are called *workloads*.

### a. Click Kubernetes Engine → Clusters → Deploy

The screenshot shows the Google Cloud console interface for Kubernetes Engine. The left sidebar has 'Kubernetes Engine' selected. The main area shows 'Kubernetes clusters' with a table listing one cluster: 'autopilot-cluster-100' in the 'us-central1' location, using 'Autopilot' mode, with 0 nodes, 0 vCPUs, and 0 GB memory. The cluster status is 'Running' (indicated by a green checkmark).

Status	Name	Location	Mode	Number of nodes	Total vCPUs	Total memory
<input checked="" type="checkbox"/>	autopilot-cluster-100	us-central1	Autopilot	0	0	0 GB

### b. In Image Path, click → Select → Artifact Registry

→ Select the hello-repo repository

→ Select hello-app image

→ Click SELECT



Select container image ✕

ARTIFACT REGISTRY CONTAINER REGISTRY

Project: midyear-nebula-406522 [CHANGE](#)

▼ us-west1-docker.pkg.dev/midyear-nebula-406522/hello-repo

▼ hello-app

aa0facf1cb v1 17 minutes ago

[SELECT](#) [CANCEL](#)

c. Specify a deployment name,

d. Refer to the YAML deployment configuration file, what are the minimum and maximum replicas values?

[←](#) Create a deployment

☒ Container

2 Configuration

A deployment is a configuration which defines how Kubernetes deploys, manages and scales your container image. Kubernetes will ensure that your system matches this configuration. Three replicas will be created by default.

Deployment name \*  
hello1

Namespace \*  
default

Labels

Use Kubernetes labels to control how workloads are scheduled to your nodes. Labels are applied to all nodes in this node pool.

Key 1 \*  
app

Value 1  
hello1

[+ ADD KUBERNETES LABEL](#)

Configuration YAML

Kubernetes deployments are defined declaratively using YAML files. The first

e. ➔ Click Deploy.



←

Replica set de...

REFRESH

EDIT

DELETE

ACTIONS ▾

KUBECTL ▾

OPERATIONS

Pods

3 current/3 desired

Label selector

app = hello1

pod-template-hash = 58ff7c7956

Managing deployment

[hello1](#)

Pod specification

Labels

app: hello1

pod-template-hash: 58ff7c7956

Termination grace period

30

Restart policy

Always

Containers

[hello-app-sha256-1](#)

Managed pods

Name	Status	Restarts	Created on ↑
<a href="#">hello1-58ff7c7956-hwkqw</a>	✓ Running	0	29 Nov 2023, 02:54:02
<a href="#">hello1-58ff7c7956-lzrh9</a>	✓ Running	0	29 Nov 2023, 02:54:02
<a href="#">hello1-58ff7c7956-bmnpv</a>	✓ Running	0	29 Nov 2023, 02:54:02

f. Click Workloads → Expose