

HW-4: Running MNIST Training & Inference on K8s

1. INTRODUCTION

In modern computing, the emergence of virtualization and containers have changed how software is developed, deployed and maintained. As computing needs of quickly growing softwares increases, virtualization and containers are heavily used for efficient use of resources and providing flexibility in deployment. However, managing and orchestrating these containers at scale presents a new set of challenges. This report created as part of *Homework-4* talks about running MNIST digit recognition model training and inference on a K8 cluster while making use of persistent storage.

2. BACKGROUND STUDY

What is Kubernetes?

Kubernetes, (K8s), was initially created to streamline the deployment and management of containerized applications. It was developed by Google and is currently an open source platform. It was based on Google's internal container orchestration system Borg. The main objectives of Kubernetes is:

- Automation: Kubernetes automates the deployment and scaling of containerized applications. This means that we as developers don't need to manually handle things and this also helps reduce human errors.
- Scalability: It does automatic scaling. This allows the applications to handle varying loads dynamically.
- Portability: Kubernetes provides a consistent environment across different infrastructure providers, making it easier to migrate and run applications seamlessly in various environments. (*I experienced this personally while working on the project as running container meant that my teammate who uses a Macbook and me who uses a Windows laptop didn't run into infrastructure issues, which we would usually run into before using containers and Kubernetes*)

Some important Kubernetes concepts and terms:

- Pods are the smallest units which can be deployed on Kubernetes. They consist of one or more than one container that share the same network and storage.
- Nodes are physical or virtual machines in a Kubernetes cluster where containers are being deployed.
- Clusters are a group of nodes.
- Services provide a stable endpoint to access a set of pods which means that we can communicate easily inside a cluster.
- Controllers handle deployment, scaling and updating of application components.

In a nutshell, containerized apps run inside pods, which are the smallest deployable units that run on nodes. Nodes form a cluster where compute, storage, and network resources are combined, shared, and coordinated across pods. Pods are grouped together in a deployment and presented as a service to clients or other applications and services

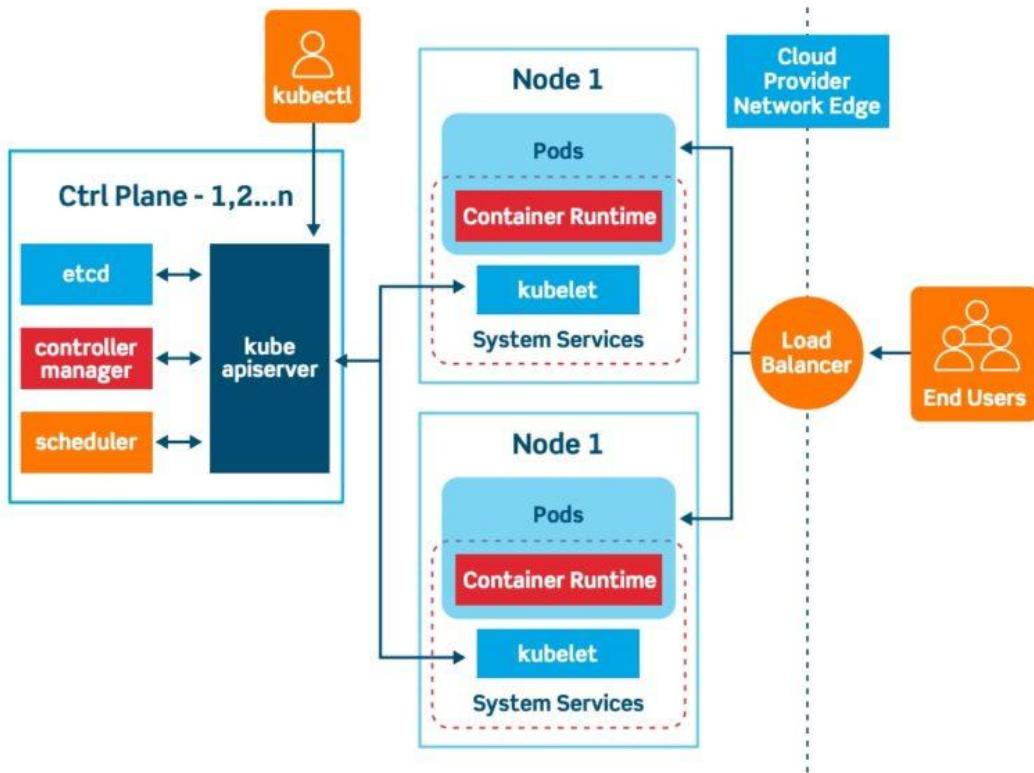


Figure: Kubernetes Architecture

GKE AUTOPILOT V/s STANDARD CLUSTERS

There are two types of clusters available for creation in GCP GKE:

- 1) Autopilot which is the automatic and hands off version
- 2) Standard Cluster which requires configuration but offers more control.

In an autopilot cluster the zonal control plane and nodes are managed and provisioned by GKE which means it automatically scales depending on demand. Standard clusters offer more control where the user is responsible for scaling

Autopilot provides a good starting point to quickly deploy applications without spending too much time understanding demands and nodes. It does have limitations with respect to the number of pods per cluster, usage of GKE managed namespace etc. but if someone wants to

quickly setup autopilot is the way to go. Standard clusters offer more control but require deep understanding of your applications demands, traffic and is often the work of a skilled developer, while autopilot is easy to use even for a beginner.

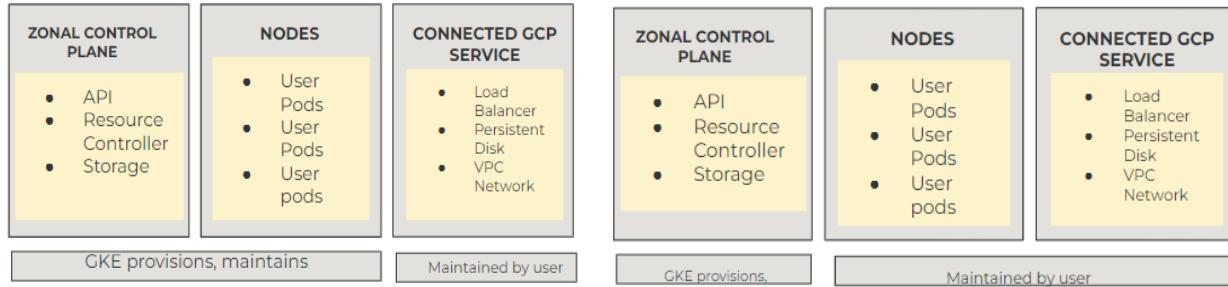


Figure: GKE Autopilot (left) , GKE Autopilot(Right)

Persistent Storage

Persistent storage is the long-term memory of the application and is required when data needs to survive beyond the lifespan of individual containers. Some reasons why persistent storage is needed are:

1. Data Retention: In many applications, data needs to persist even if the application or container is restarted or moved to a different node. Without persistent storage, all that valuable data would be lost each time a container shuts down.
2. Stateful Applications: Some applications have a state that needs to be preserved. If we need to ensure that your application can recover its state after being restarted or moved, we need persistent storage.
3. Scaling Applications: When applications need to be scaled horizontally, persistent storage ensures that each instance of the application has the same data. This ensures that data consistency is maintained to avoid any loss of data integrity.

Key Persistent Storage Concepts wrt Google Kubernetes Engine:

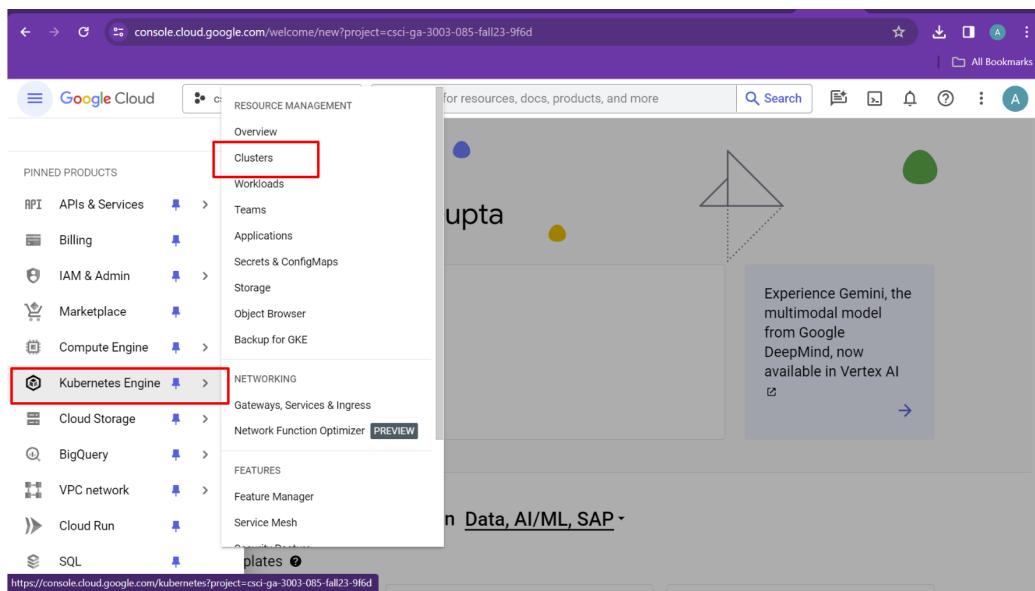
1. Persistent Volumes (PVs): PVs are the abstraction layer that represents physical storage resources in a cluster. They decouple the details of how storage is provided from how it's consumed. A PV in Kubernetes could be a physical disk or even a network-attached storage.
2. Persistent Volume Claims (PVCs): PVCs act as a request for storage by a user or a group of users. They provide a way to claim a specific amount of storage defined by the PV.

3. Storage Classes: These are used to dynamically provision PVs based on predefined templates. Storage Classes enable automatic provisioning, ensuring that applications get the storage they need without manual intervention.

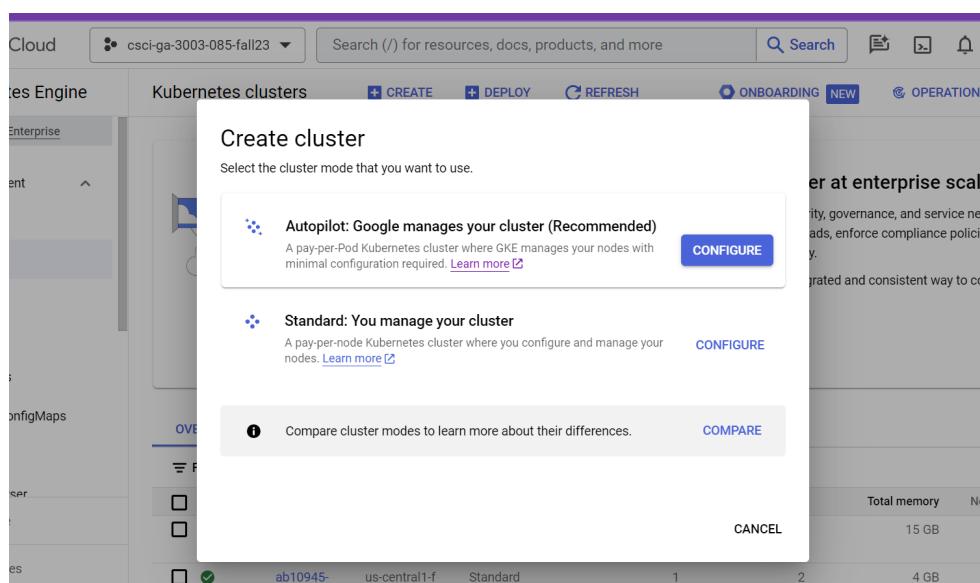
3. CREATING A KUBERNETES CLUSTER ON GCP

I created a cluster using Autopilot using the GCP user interface.

1) Select Clusters



2) Create an autopilot cluster



3) Select cluster settings.

The screenshot shows the 'Create an Autopilot cluster' wizard on the Google Cloud Platform. The left sidebar lists five steps: 1. Cluster basics (selected), 2. Networking, 3. Advanced settings, 4. Fleet, and 5. Review and create. The main panel is titled 'Cluster basics' and contains the following information:

- Create an Autopilot cluster by specifying a name and region. After the cluster is 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-2 (input field)
- Region:** us-central1 (dropdown menu)

At the bottom are buttons for 'NEXT: NETWORKING', 'RESET SETTINGS', 'CREATE' (highlighted in blue), 'CANCEL', and links to 'REST' or 'COMMAND LINE'.

Using command line-

```
gcloud beta container --project "csci-ga-3003-085-fall23-9f6d" clusters
create-auto "ag8733-cluster1" --region "us-central1" --release-channel
"regular" --network
"projects/csci-ga-3003-085-fall23-9f6d/global/networks/csci-ga-3003-085-fall23-
net" --subnetwork
"projects/csci-ga-3003-085-fall23-9f6d/regions/us-central1/subnetworks/csci-ga-
3003-085-fall23-subnet-02" --cluster-ipv4-cidr "/17"
--binauthz-evaluation-mode=DISABLED
```

I created an autopilot cluster called ag8733-cluster1, which automatically scales the nodes up and down.

Status	Name	Location	Mode	Number of nodes	Total vCPUs	Total memory	Notifications
<input type="checkbox"/>	ab10945-hw4	us-central1-f	Standard	1	4	15 GB	-
<input type="checkbox"/>	ab10945-hw4-cluster2	us-central1-f	Standard	1	2	4 GB	-
<input type="checkbox"/>	ag8733-cluster1	us-central1	Autopilot	1	4 GB	-	-
<input type="checkbox"/>	am13018-cluster-1	us-central1	Autopilot	1	4 GB	-	-
<input type="checkbox"/>	autopilot-cluster-1	us-central1	Autopilot	0	0 GB	-	-

Figure: Created cluster ag8733-cluster1

4. SETTING UP PVC

```

apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: ag8733pvc
spec:
  accessModes:
    - ReadWriteOnce
  resources:
    requests:
      storage: 7Gi
  storageClassName: standard-rwo

```

Figure: pvc.yaml

- apiVersion: SPecifies the API version of the resource, here I am using the core/v1 version of the Kubernetes API
- Kind : this defines the type of resource. Here, it's a PersistentVolumeClaim.
- The name of the PVC has been set using the metadata field. I have set it as ag8733pvc.
- accessModes: Describes the desired access mode for the persistent volume. ReadWriteOnce, means the volume can be mounted as read-write by a single node.
- Resources are set to 7Gi meaning that I'm requesting for a 7GB storage.

- `storageClassName`: set to “standard-rwo.”

The screenshot shows the Google Cloud Kubernetes Engine Storage interface. On the left, there's a sidebar with options like Clusters, Workloads, Teams, Applications, Secrets & ConfigMaps, Storage (which is selected), Object Browser, and Backup for GKE. The main area has tabs for Storage, Refresh, and Delete. It shows a cluster dropdown set to "cloudml-gke-cluster-1, ag...", a Namespace dropdown, and buttons for RESET and SAVE. Below these are two tabs: PERSISTENT VOLUME CLAIMS (selected) and STORAGE CLASSES. A message says "Persistent volume claims are requests for storage of specific size and access mode." Below this is a table with a red border around its header and rows. The columns are Name (sorted by name), Phase, Volume, Storage class, Namespace, and Cluster. There are three entries:

Name	Phase	Volume	Storage class	Namespace	Cluster
ag8733-pvc	Bound	pvc-86eb47ac-29f9-4d9d-a896-ada6538390d3	standard-rwo	default	ag8733-cluster1
am13018-pvc	Bound	pvc-11b40d72-f276-4632-8168-80a603bc6c8e	standard-rwo	default	am13018-cluster-1

Figure: Shows my persistent storage which has been bound.

5. TRAINING THE MODEL

The main goal of this exercise was to create a containerized version of the training script which can do the training for the MNIST recognizer model and then save the weights of the model to persistent storage so that they can be later used during inference.

- 1) I made the necessary changes in the pytorch/examples mnist.py file.
I train the model for one epoch and then save the model weights in the `/mnt` directory with the name `ag8733_model.pth`
- 2) I then create a dockerfile which uses the pytorch base image and then runs the `train.py` script. I use this command inside the dockerfile to run the script-

```
CMD python train.py --batch-size 32 --epochs 1 --no-cuda  
--no-mps --save-model
```

- 3) Build the docker file using docker build command and then tag the image.
- 4) Create a repository on dockerhub
- 5) Push this train image to docker hub using-

```
Docker push anoushkgupta/mnisttrain:v3
```

```

dataset1 = datasets.MNIST('../data', train=True, download=True,
                        transform=transform)
dataset2 = datasets.MNIST('../data', train=False,
                        transform=transform)
train_loader = torch.utils.data.DataLoader(dataset1, **train_kwargs)
test_loader = torch.utils.data.DataLoader(dataset2, **test_kwargs)

model = Net().to(device)
optimizer = optim.Adadelta(model.parameters(), lr=args.lr)

scheduler = StepLR(optimizer, step_size=1, gamma=args.gamma)
for epoch in range(1, args.epochs + 1):
    train(args, model, device, train_loader, optimizer, epoch)
    #test(model, device, test_loader)
    scheduler.step()

    if args.save_model:
        print("Saving model!")
        torch.save(model.state_dict(), "/mnt/ag8733_model.pth")

```

Figure: changes made to the train.py file to just perform training and save weights

RUNNING THE TRAINING JOB ON GCP:

1. I used train.yaml to start a job which makes use of the persistent storage which was created earlier.
2. Use command **kubectl -f apply train.yaml** to run, which creates the job in the K8 cluster.
3. This job uses the image I had pushed to dockerhub.
4. Once this job finishes, the container dies but as weights are saved in persistent storage they persist.

Severity	Timestamp	SUMMARY
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [49600/60000 (83%)] Loss: 0.063072
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [49920/60000 (83%)] Loss: 0.073397
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [50240/60000 (84%)] Loss: 0.108103
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [50560/60000 (84%)] Loss: 0.066630
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [50880/60000 (85%)] Loss: 0.118702
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [51200/60000 (85%)] Loss: 0.048884
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [51520/60000 (86%)] Loss: 0.218634
> i	2023-12-20 18:57:47.276 EST	Train Epoch: 1 [51840/60000 (86%)] Loss: 0.010028

The screenshot shows the Kubernetes UI for a pod named 'ag8733trainjob5'. The 'LOGS' tab is selected. The logs output shows the following training progress:

```

Kubernetes node audit logs from the control plane
Severity: Default
Container logs
SUMMARY
Severity | Timestamp | Train Epoch: 1 [57920/60000 (97%)] | Loss: 0.004991
Severity | Timestamp | Train Epoch: 1 [58240/60000 (97%)] | Loss: 0.000312
Severity | Timestamp | Train Epoch: 1 [58560/60000 (98%)] | Loss: 0.002836
Severity | Timestamp | Train Epoch: 1 [58880/60000 (98%)] | Loss: 0.000497
Severity | Timestamp | Train Epoch: 1 [59200/60000 (99%)] | Loss: 0.000534
Severity | Timestamp | Train Epoch: 1 [59520/60000 (99%)] | Loss: 0.000301
Severity | Timestamp | Train Epoch: 1 [59840/60000 (100%)] | Loss: 0.004318
Severity | Timestamp | Saving model!

```

Figure: Training Logs

```

ag8733@cloudshell:~ (csci-ga-3003-085-fall23-9f6d)$ kubectl get pods
NAME                               READY   STATUS    RESTARTS   AGE
ag8733inferencejob-5f4bdbcb9-vthl8  1/1    Running   0          101m
ag8733inferencejob3-6fb779cb78-ttrzn 1/1    Running   0          116m
aq8733testjob4-rxlzf                0/1    Completed  0          4h3m
ag8733trainjob5-7jlqt               0/1    Completed  0          4h14m
ag8733@cloudshell:~ (csci-ga-3003-085-fall23-9f6d)$

```

Figure: kubectl get pods shows the train job has been completed

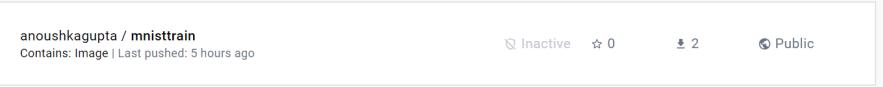


Figure: The Dockerhub repository for the training image

6. INFERENCE

For inference I used the flask app I created for the Final project and use the same templates for deploying the inference service of HW-4.

The Flask app uses the weights stored in persistent storage to carry out inference and this is provisioned using the infer.yaml file.

I then create an endpoint by using service.yaml.

```

File Edit Selection View Go Run
! pvc.yaml train.py Dockerfile app.py X
C: > Users > Anoushka Gupta > Documents > NYU > Fall2023 > Cloud > hw4 > flask_app > app.py
48     @app.route('/predict', methods=['POST'])
49     def predict():
50         res = {}
51         file = request.files['file']
52         print("here 1.0")
53         if not file:
54             print("here 1")
55             res['message'] = 'Error. Cannot find image'
56             res['Result'] = 'NA'
57         else:
58             print("here 2")
59             res['message'] = 'Success'
60             image = Image.open(file)
61             image = transform(image).unsqueeze(0)
62             res['Result'] = predict2(image)
63             print("Returning result json")
64     return res['Result']
65 model = Net()
66 model.load_state_dict(torch.load("/mnt/ag8733_model.pth"))
67 model.eval()
68

```

Figure: Flask app which uses the weights in the persistent storage

The Dockerfile just runs this flask application.

The inference.yaml runs this on a pod. The name of the job is ag8733inferencejob.

```

ag8733@cloudshell:~ (csci-ga-3003-085-fall23-9f6d)$ kubectl get pods
NAME          READY   STATUS    RESTARTS   AGE
ag8733inferencejob-5f4bdb89-vthl8   1/1     Running   0          109m
ag8733inferencejob3-6fb779cb78-dr86b 0/1     Pending    0          14s
ag8733inferencejob3-6fb779cb78-ttrzn 1/1     Terminating   0          123m
ag8733testjob4-rxlzf                0/1     Completed   0          4h11m
ag8733trainjob5-7jlqt               0/1     Completed   0          4h21m
ag8733@cloudshell:~ (csci-ga-3003-085-fall23-9f6d)$

```

I then use service.yaml to create an end point which listens on port 8000.

The service.yaml is run using **kubectl -f apply service.yaml** command.



Figure: The DockerHub repository for the inference image

The screenshot shows the Kubernetes Engine dashboard for a deployment named 'ag8733-cluster1'. The deployment has 0 replicas. It is associated with the 'default' namespace, has labels 'app: serviceapi' and 'tier: backend', and uses 'Container logs' and 'Audit logs' for monitoring. The deployment is currently at revision 1, which contains one container named 'ag8733inference' and a volume 'ag8733pvc'. The deployment is not configured for an autoscaler or vertical pod autoscaler.

Active revisions	Revision	Name	Status	Summary	Created on	Pods running/Pods total
	1	ag8733inferencejob-5f4bdbc89	OK	ag8733inference: anoushkagupta/mnistinfer:v4	Dec 20, 2023, 9:24:22 PM	1/1

Figure: Dashboard of the inference job, We can see the volume that is mounted the no of pods running, namespace, cluster etc.

The screenshot shows the 'Exposing services' section of the Kubernetes Engine dashboard. It lists three services: 'ag8733inferencejob-cmw8s' (Load balancer, endpoints 34.134.123.20:80), 'ag8733inferencejob3' (Load balancer, endpoints 104.154.104.171:80), and 'ag8733inferencejob3-4rlsur' (Load balancer, endpoints 34.66.123.100:80).

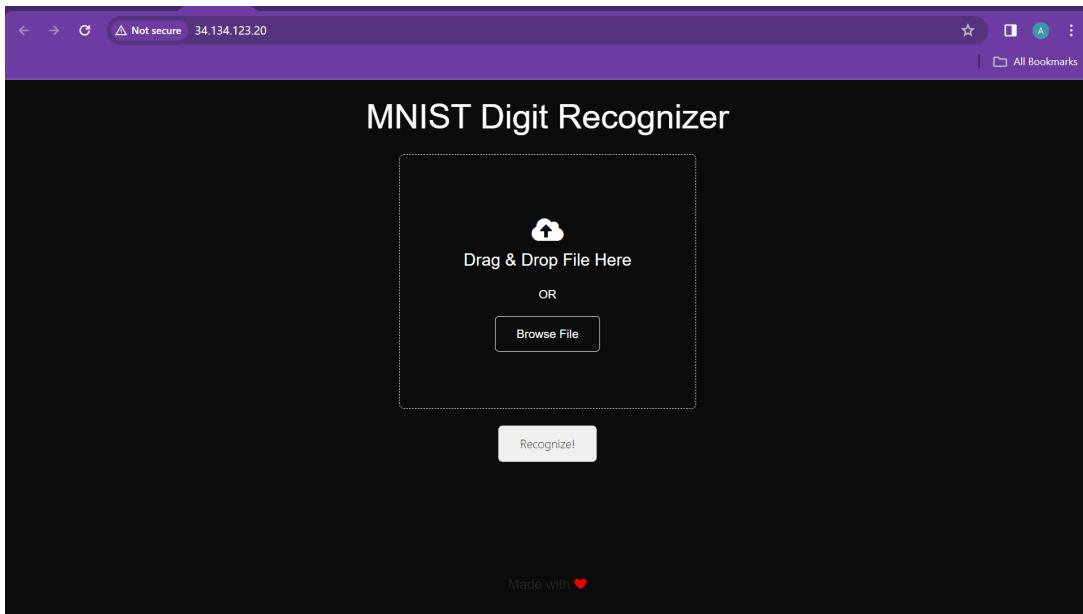
Name	Type	Endpoints
ag8733inferencejob-cmw8s	Load balancer	34.134.123.20:80
ag8733inferencejob3	Load balancer	104.154.104.171:80
ag8733inferencejob3-4rlsur	Load balancer	34.66.123.100:80

Figure: The exposed services with this inference pods

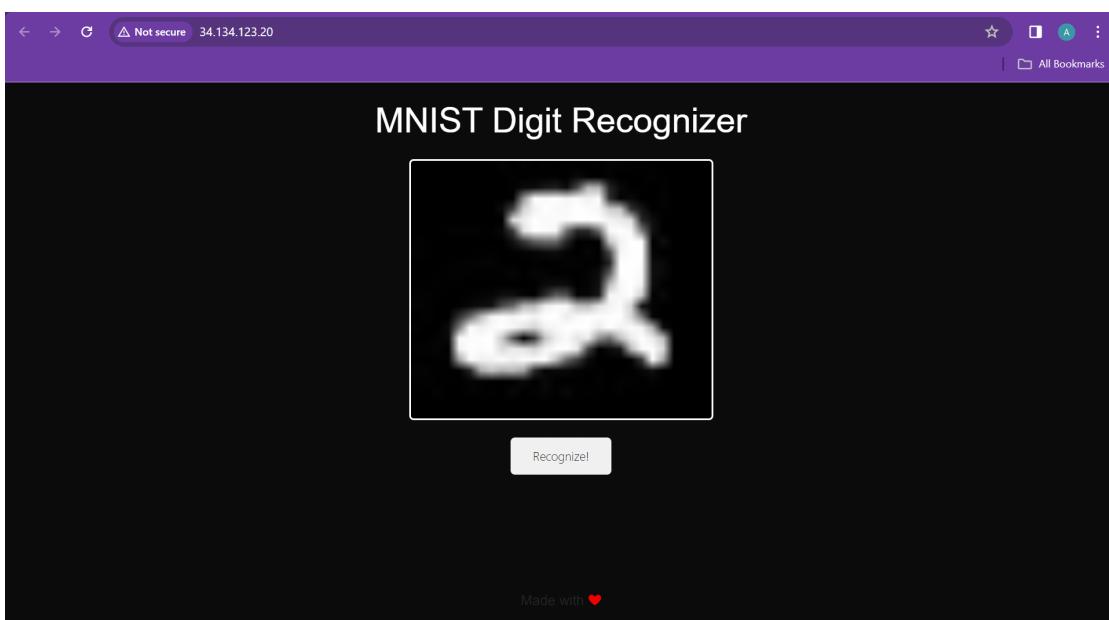
7. DEMO

The end point of the service is <http://34.134.123.20/> where you can upload a file and it shows the predicted digit.

UI of the flask app-



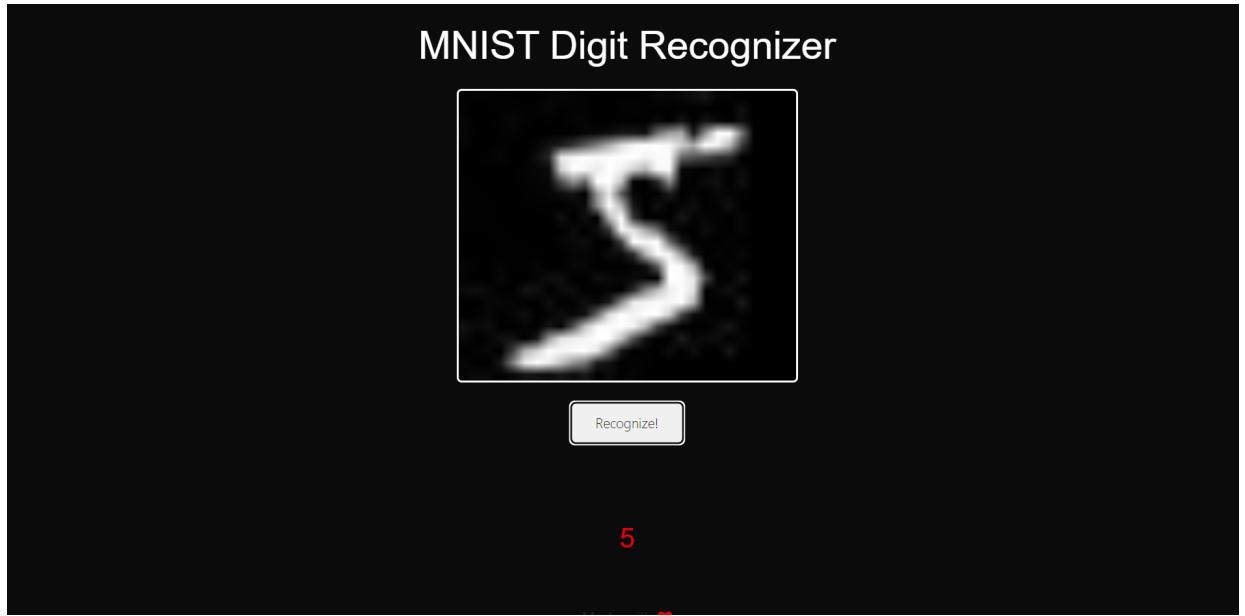
Uploading a MNIST dataset image- Click on the button Recognize to generate results.



The predicted result can be seen on the screen at the bottom.



Trying with another example-



8. CONCLUSION & LEARNINGS

This was a very good deep down into training and saving weights in persistent storage. Some challenges that I encountered included-

- 1) Setting up a persistent storage and ensuring the weights are actually being stored in the mounted volume. I wasn't specifying the correct paths while storing the model weights and hence they were not being saved on the persistent volume and the weights would die when the training container ended.
- 2) Requesting Pods. There was a lot of traffic and it was taking a very long time before pods were being assigned to my autopilot cluster and this significantly slowed down my progress.

Learnings-

- 1) Setting up machine learning training jobs using containers and Kubernetes on GCP.
- 2) Command line understanding of kubectl and associated commands like apply, get pods, get service etc.

In conclusion I did the following things-

- Created a persistent storage volume claim using pvc.yaml
- Created a container for the training part of code, pushed the image to docker hub, wrote the train.yaml script and then ran a job for training which saves model weights.
- Created a flask app which uses the weights stored in persistent storage for inference.
- Created a deployment or inference.yaml to run the flask application on a pod.
- Expose the endpoint for communication using service.yaml.

REFERENCES

1. Understanding Kubernetes concepts
(<https://platform9.com/blog/kubernetes-enterprise-chapter-2-kubernetes-architecture-concepts/>)
2. Docker Documentation (<https://docs.docker.com/>)
3. Deploying an app on a GKE cluster
(<https://cloud.google.com/kubernetes-engine/docs/deploy-app-cluster>)
4. Pytorch MNIST example (<https://github.com/pytorch/examples/tree/main/mnist>)
5. Persistent Volume on GCP
(<https://cloud.google.com/kubernetes-engine/docs/concepts/persistent-volumes>)
6. YAML files for Kubernetes
(<https://www.cncf.io/blog/2022/03/03/how-to-write-yaml-file-for-kubernetes/>)