

Automated Deep-Learning based Image Captioning System with Dev-Ops Driven Deployment

**Software Production Engineering - Final Project
IIITB M.Tech 2nd Sem, 2025**

**Aman Bahuguna
(MT2024018)**

**Soumik Pal
(MT2024153)**

DevOps is invaluable in today's IT topology, because it bridges the gap between software development and IT operations, enabling faster and more reliable delivery of applications and services. By fostering collaboration, automating repetitive tasks like testing, deployment, and infrastructure management, DevOps reduces errors and accelerates release cycles. This leads to improved product quality, quicker feedback loops, and greater agility in responding to changing business needs, ultimately driving better customer satisfaction and competitive advantage.

In this project, we design and implement a comprehensive DevOps framework to fully automate the Software Development Life Cycle (SDLC)

using cutting-edge DevOps tools. Focusing on the specialized domain of Machine Learning, we seamlessly integrate DevOps principles to unlock its benefits—thus pioneering a robust and efficient MLOps pipeline.

1. ML Component :

Problem Statement :

Given an image, generate a caption for it which would serve as the textual description of the image. This task combines techniques from both Computer Vision and Natural Language Processing, making it a multidisciplinary challenge.

Model Architecture chosen :

Typically, image captioning models follow an encoder-decoder architecture—where the image is first transformed into a meaningful feature representation by the encoder, and then the decoder translates this representation into a coherent descriptive sentence.

In our project, we have used a DenseNet CNN to extract image features and an LSTM to generate text by combining image and word embeddings.

Dataset Used : <https://www.kaggle.com/datasets/adityajn105/flickr8k>

We used a subset of the standard Flickr dataset, called Flickr8K that has ~8K images and 5 captions for each image, creating around 40K image, caption pairs. This was split into a training dataset, validation dataset and testing dataset (80% - 10% - 10% split).

Preprocessing :

To make the training, validation and testing process quicker, the feature extraction of all 8k images, using DenseNet201 model from tensorflow library, were done beforehand and stored as .pkl artifact in models folder. This was later loaded and used during training, validation and testing as a simple lookup.

The caption text preprocessing involved cleaning and normalizing sentences by lowercasing, removing special characters, extra spaces, and single letters, and adding start and end tags for model training.

Training the Model :

The model was trained over 12 epochs using a carefully selected learning rate and batch size to optimize performance. To ensure efficient training and prevent overfitting, model checkpointing was employed to save the best-performing model based on validation loss, allowing for recovery and reuse of the most effective weights. Additionally, early stopping was implemented with a patience of 3, meaning training would halt if the validation loss did not improve for three consecutive epochs—thereby saving time and computational resources while maintaining model generalization.

The model and tokenizer used were saved in the local models directory to be later reused during testing and inference.

Testing the Model :

The save model was loaded and tested on the test dataset to evaluate BLEU score. The BLEU score evaluates the quality of generated captions by measuring their similarity to reference captions based on overlapping n-grams.

Inference and Frontend App :

The inference script has no main method and instead exposes its method as an API which is called from [app.py](#), a streamlit based frontend where a user can submit their image and the api call will return the generated caption which is then displayed to the user.

Project Structure based on ML Components :

Scripts :

- `extract_features.py`
- [train.py](#)
- [test.py](#)
- [infer.py](#)

Dataset :

- `data/flickr/captions.txt`
- `data/flickr/images`

Model and Artifacts :

- `Features.pkl`
- `Tokenizer.pkl`
- `model.h5`

```
(ml-devops-env) soumik@soumik-VirtualBox:~/SPE_Final_Project$ ls extract_features.py train.py test.py infer.py app.py data/ models/
app.py  extract_features.py  infer.py  test.py  train.py

data/:
flickr

models/:
features.pkl  model.h5  tokenizer.pkl
```

The entire Machine Learning project was developed inside a Virtual Environment to ensure all dependencies could be replicated consistently.

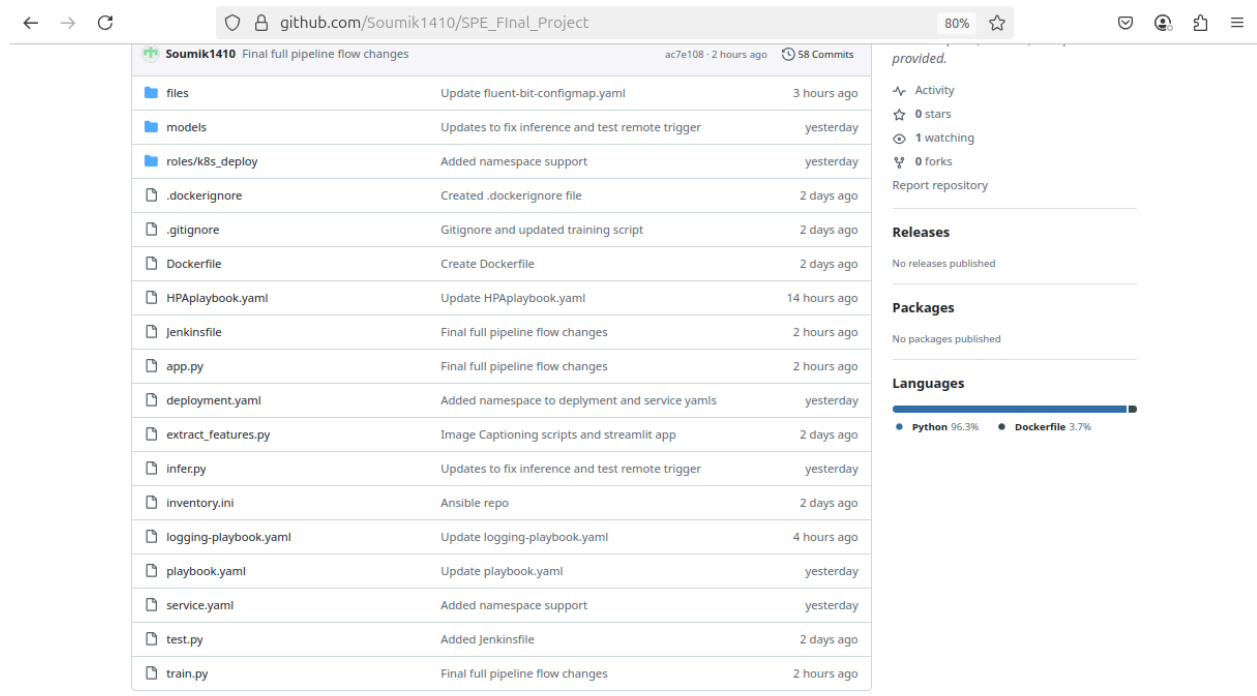
2. Code versioning - Git, Github

Git is a distributed version control system for tracking code changes, while GitHub is a cloud-based platform for hosting and collaborating on Git repositories.

In our project we are versioning all the scripts used in the ML pipeline such as `extract_features.py`, [train.py](#), [test.py](#), [infer.py](#) & [app.py](#) along with the certain artifacts that are saved during training which will be later needed during testing and inference such as the tokenizer fit on training data captions for word to embedding mappings.

Along with these, other scripts used later in the pipeline such as Jenkinsfile, Dockerfile, Ansible playbooks, inventory file, Ansible Roles, Kubernetes Manifests to deploy the application, enable logging and monitoring and HPA are also versioned in this repository.

Final project repository :



← → ↻ github.com/Soumik1410/SPE_Final_Project 80% ☆

Soumik1410 Final full pipeline flow changes ac7e108 · 2 hours ago 58 Commits

File	Commit Message	Time Ago
files	Update fluent-bit-configmap.yaml	3 hours ago
models	Updates to fix inference and test remote trigger	yesterday
roles/k8s_deploy	Added namespace support	yesterday
.dockerignore	Created .dockerignore file	2 days ago
.gitignore	Gitignore and updated training script	2 days ago
Dockerfile	Create Dockerfile	2 days ago
HPAplaybook.yaml	Update HPAplaybook.yaml	14 hours ago
Jenkinsfile	Final full pipeline flow changes	2 hours ago
app.py	Final full pipeline flow changes	2 hours ago
deployment.yaml	Added namespace to deployment and service yamls	yesterday
extract_features.py	Image Captioning scripts and streamlit app	2 days ago
infer.py	Updates to fix inference and test remote trigger	yesterday
inventory.ini	Ansible repo	2 days ago
logging-playbook.yaml	Update logging-playbook.yaml	4 hours ago
playbook.yaml	Update playbook.yaml	yesterday
service.yaml	Added namespace support	yesterday
test.py	Added Jenkinsfile	2 days ago
train.py	Final full pipeline flow changes	2 hours ago

provided.

Activity

0 stars

1 watching

0 forks

Report repository

Releases

No releases published

Packages

No packages published

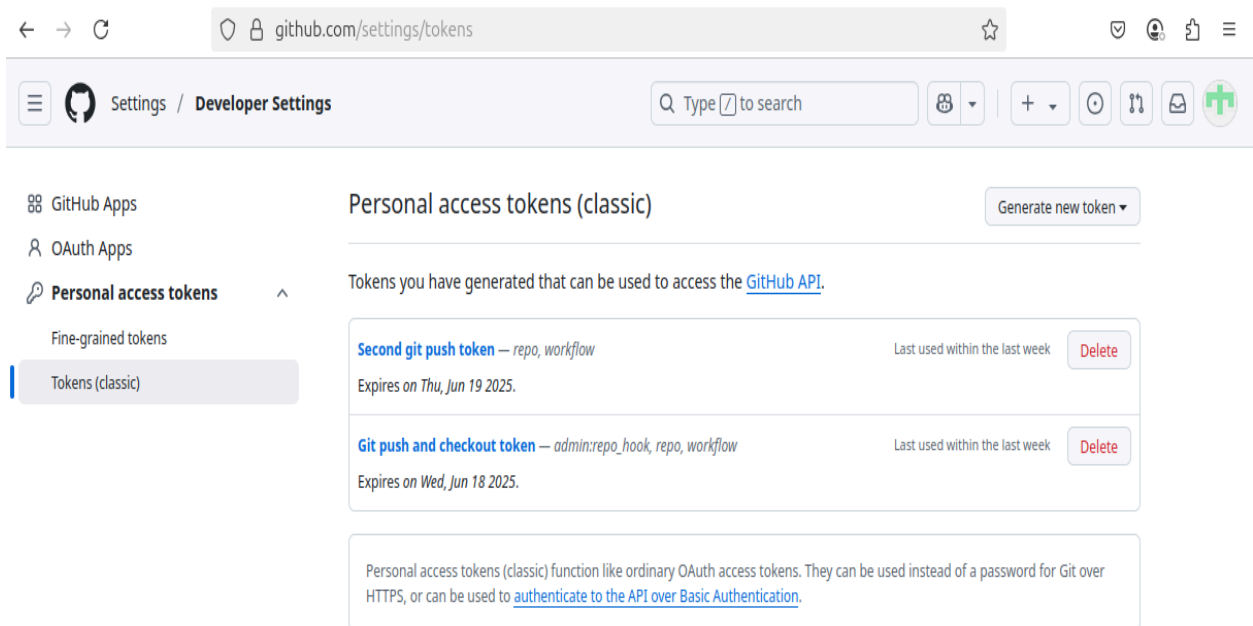
Languages

Python 96.3% Dockerfile 3.7%

Other large files present in the local development virtual environment such as the dataset and the trained models were not pushed to the Github repository by maintaining a .gitignore file

```
(ml-devops-env) soumik@soumik-VirtualBox:~/SPE_Final_Project$ cat .gitignore
__pycache__/  
*.pyc  
*.log  
*.h5  
data/flickr/Images/  
data/flickr/captions.txt  
mlruns/  
.env  
  
(ml-devops-env) soumik@soumik-VirtualBox:~/SPE_Final_Project$
```

Personal Access Tokens were created in Github for pushing commits from the local repository to the remote Github repository.



3. MLFlow - Experiment Tracking and Hyperparameter Versioning

When experimenting with different model architectures and different hyperparameters, it helps to track every experiment, the architecture & the different hyperparameters used for training, the metrics evaluated etc. so that the same model and results can be reproduced consistently and different metrics of different experiments can be easily compared.

MLFlow provides an easy way to do this, saving all these to a central server and provides a UI to view these experiments and corresponding saved parameters & metrics later. MLFlow also provides a registry to save and version models, but we opted to do this manually with checkpointing.

To efficiently track training experiments, the training loop was encapsulated within a code block that initiated an MLflow run each time a model was trained. During this process, key parameters such as batch size, optimizer, embedding dimension, LSTM units, dropout rate, and the number of training epochs were logged. At the end of each epoch, both the training loss and validation loss were recorded. Once the training was completed, the training loss versus validation loss curve was saved as an MLflow artifact, providing a visual representation of the model's performance throughout the training process.

```
# MLflow logging
with mlflow.start_run():
    # Log hyperparameters
    mlflow.log_param("batch_size", 64)
    mlflow.log_param("optimizer", "adam")
    mlflow.log_param("embedding_dim", 256)
    mlflow.log_param("lstm_units", 256)
    mlflow.log_param("dropout", 0.5)
    mlflow.log_param("epochs", 10)

    # Train and capture history
    history = caption_model.fit(
        train_gen,
        validation_data=val_gen,
        epochs=10,
        callbacks=[checkpoint, earlystop, reduce_lr],
        verbose=1
    )

    # Log metrics per epoch
    for epoch in range(len(history.history['loss'])):
        mlflow.log_metric("loss", history.history['loss'][epoch], step=epoch)
        mlflow.log_metric("val_loss", history.history['val_loss'][epoch], step=epoch)
```

```

# Log metrics per epoch
for epoch in range(len(history.history['loss'])):
    mlflow.log_metric("loss", history.history['loss'][epoch], step=epoch)
    mlflow.log_metric("val_loss", history.history['val_loss'][epoch], step=epoch)

# Save training curve
plt.figure()
plt.plot(history.history['loss'], label='train_loss')
plt.plot(history.history['val_loss'], label='val_loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.title('Training Curve')
plot_path = os.path.join(model_dir, 'training_curve.png')
plt.savefig(plot_path)
mlflow.log_artifact(plot_path)

```

Similarly, in the test script, the BLEU score evaluated on the trained model loaded in the script is logged along with the length of the test dataset & the complete set of predictions made on the test set.

To view these saved metrics and parameters, we activate the UI with the command : mlflow ui

```

(ml-devops-env) soumik@soumik-VirtualBox:~/SPE_Final_Project$ mlflow ui
[2025-05-21 16:36:43 +0530] [2847600] [INFO] Starting gunicorn 23.0.0
[2025-05-21 16:36:43 +0530] [2847600] [INFO] Listening at: http://127.0.0.1:5000 (2847600)
[2025-05-21 16:36:43 +0530] [2847600] [INFO] Using worker: sync
[2025-05-21 16:36:43 +0530] [2847625] [INFO] Booting worker with pid: 2847625
[2025-05-21 16:36:43 +0530] [2847626] [INFO] Booting worker with pid: 2847626
[2025-05-21 16:36:44 +0530] [2847627] [INFO] Booting worker with pid: 2847627
[2025-05-21 16:36:44 +0530] [2847642] [INFO] Booting worker with pid: 2847642

```

This boots up the UI on port 5000 on localhost. We can open this in a browser to view the different saved runs and their corresponding parameters and metrics.

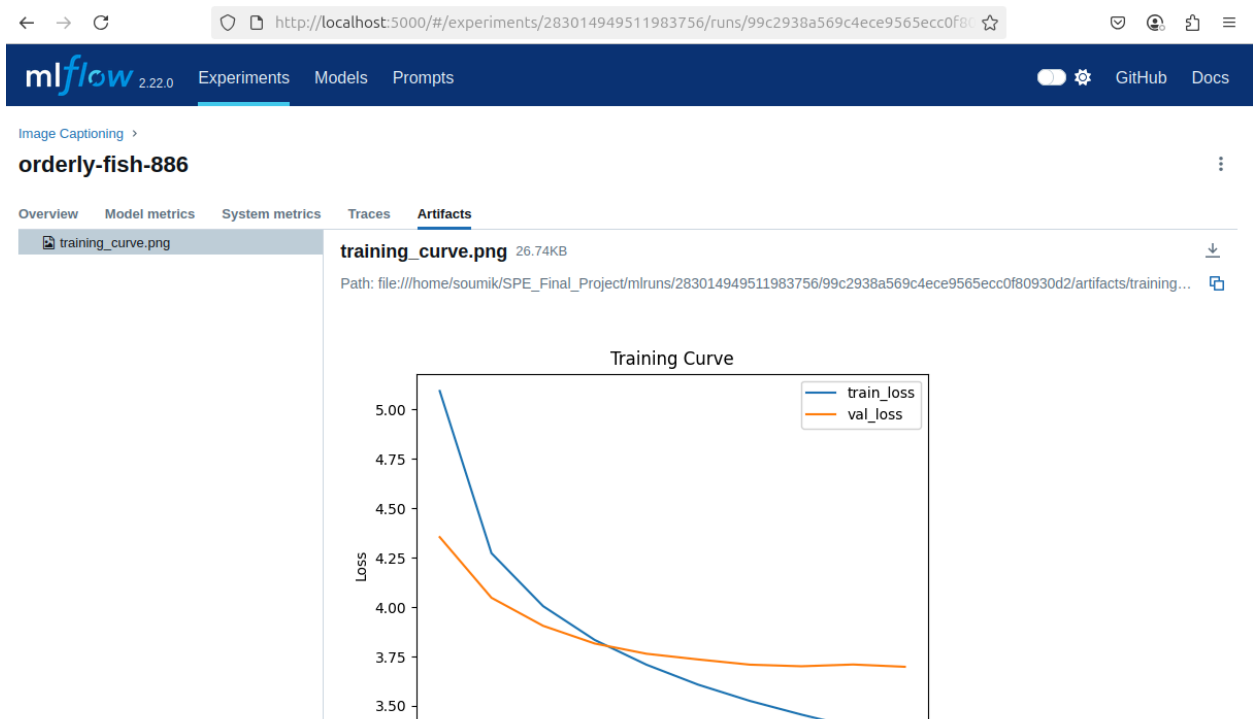
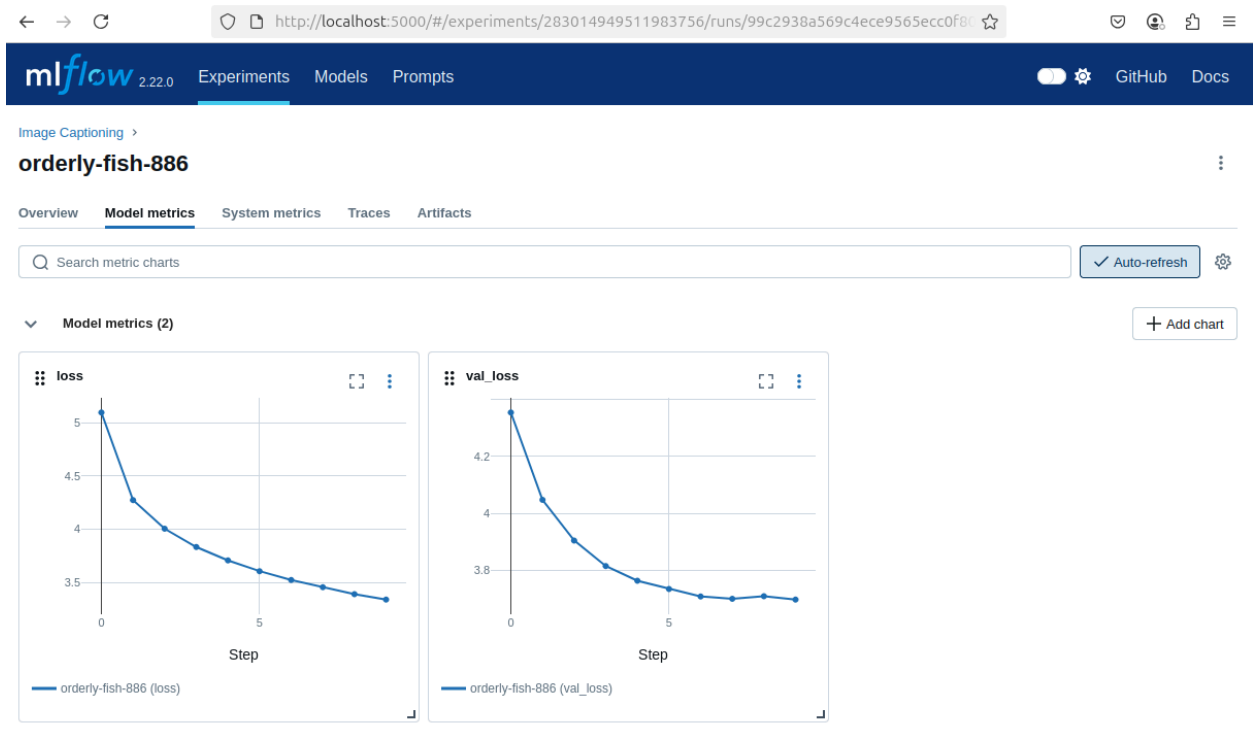
Share

14 matching runs

...

Artifacts

Metric	Value
val_loss	3.6981589794158936
loss	3.34120774269104



Experiments

Search experiments

☐ Default ✎ 🗑

☒ Image Captioning Test Eval... ✎ 🗑

☐ Image Captioning ✎ 🗑

Image Captioning Test Evaluation

[Provide Feedback](#) [Add Description](#)

Share

Runs Evaluation **Experimental** Traces

📄 📊 🔍 Time created ▾ ⋮ + New run

State: Active ▾ Datasets ▾ Sort: Created ▾ Columns ▾ Group by ▾

<input type="checkbox"/>	Run Name	Created ↕	Dataset	Duration	Source
<input type="checkbox"/>	test_evaluation	🕒 1 hour ago	-		🏠 test.py
<input type="checkbox"/>	test_evaluation	✅ 22 hours ago	-	1.3h	🏠 test.py
<input type="checkbox"/>	test_evaluation	✅ 1 day ago	-	1.4h	🏠 test.py
<input type="checkbox"/>	test_evaluation	🕒 1 day ago	-		🏠 test.py
<input type="checkbox"/>	test_evaluation	❌ 1 day ago	-	1.1h	🏠 test.py

5 matching runs

Image Captioning Test Evaluation >

test_evaluation

⋮

Overview **Model metrics** System metrics Traces Artifacts

Registered models	—
Registered prompts	—

Parameters (1)

Parameter	Value
test_samples	4050

Metrics (1)

Metric	Value
average_bleu	0.041442931007944314

The screenshot shows the mlflow web interface at localhost:5000. The 'test_evaluation' experiment is selected, and the 'Artifacts' tab is active. The 'test_predictions.csv' artifact is highlighted, and a preview of the first 500 rows is shown in a table.

image	caption	predicted_caption
436015762_8d0bae90c3.jpg	startseq man prepares to enter the red bui...	startseq man is standing on the street end...
436015762_8d0bae90c3.jpg	startseq man walking around the corner of...	startseq man is standing on the street end...
436015762_8d0bae90c3.jpg	startseq man walks past red building with ...	startseq man is standing on the street end...
436015762_8d0bae90c3.jpg	startseq man walks under building with lar...	startseq man is standing on the street end...
436015762_8d0bae90c3.jpg	startseq person walking by red building wi...	startseq man is standing on the street end...
436393371_822ee70952.jpg	startseq brown doberman is outside with s...	startseq dog is running through the grass ...

4. Jenkins - Automation CI/CD server

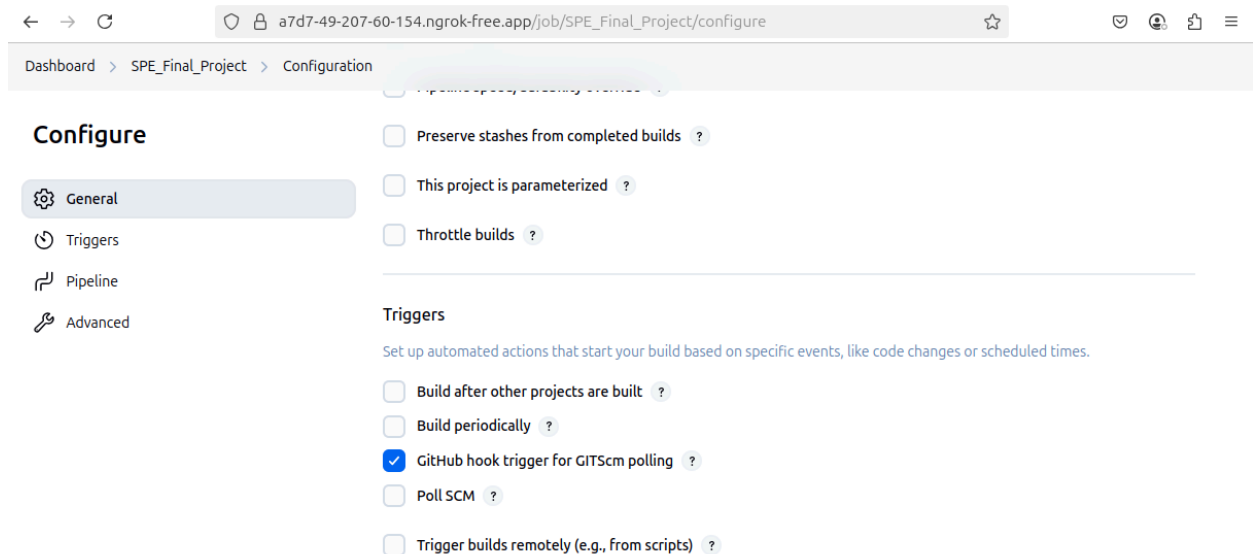
Jenkins is an open-source automation server widely used in DevOps for continuous integration and continuous delivery (CI/CD). It facilitates the automation of software development processes by allowing developers to automatically build, test, and deploy code changes.

In our project, we developed a Jenkins pipeline that is triggered by webhooks sent by a SCM service (Github), and it performs the following actions automatically :

1. Checkout the latest code from the repository
2. Train the model

3. Execute the automated test script to evaluate the trained model
4. Generate a list of the environment's dependencies
5. Prune the list to include only the essential dependencies
6. Dockerize the application
7. Push the built image to Docker Hub
8. Trigger an Ansible playbook to pull the image and deploy it on the Kubernetes cluster

For this, we will go over the different stages of the pipeline as and when we discuss the other DevOps tools used. For now, we will show how the pipeline was set up to trigger from Github webhooks and the first 3 stages till automated testing.



← → ↻ a7d7-49-207-60-154.ngrok-free.app/job/SPE_Final_Project/configure ☆ 🔒 📄 ☰

Dashboard > SPE_Final_Project > Configuration

Define your Pipeline using Groovy directly or pull it from source control.

Configure

- General
- Triggers
- Pipeline**
- Advanced

Definition

Pipeline script from SCM

SCM ?

Git

Repositories ?

Repository URL ?

https://github.com/Soumik1410/SPE_Final_Project/

Credentials ?

Soumik1410/*****

+ Add

Advanced

← → ↻ a7d7-49-207-60-154.ngrok-free.app/job/SPE_Final_Project/configure ☆ 🔒 📄 ☰

Dashboard > SPE_Final_Project > Configuration

Configure

- General
- Triggers
- Pipeline**
- Advanced

Branches to build ?

Branch Specifier (blank for 'any') ?

*/master

Add Branch

Repository browser ?

(Auto)

Additional Behaviours

Add

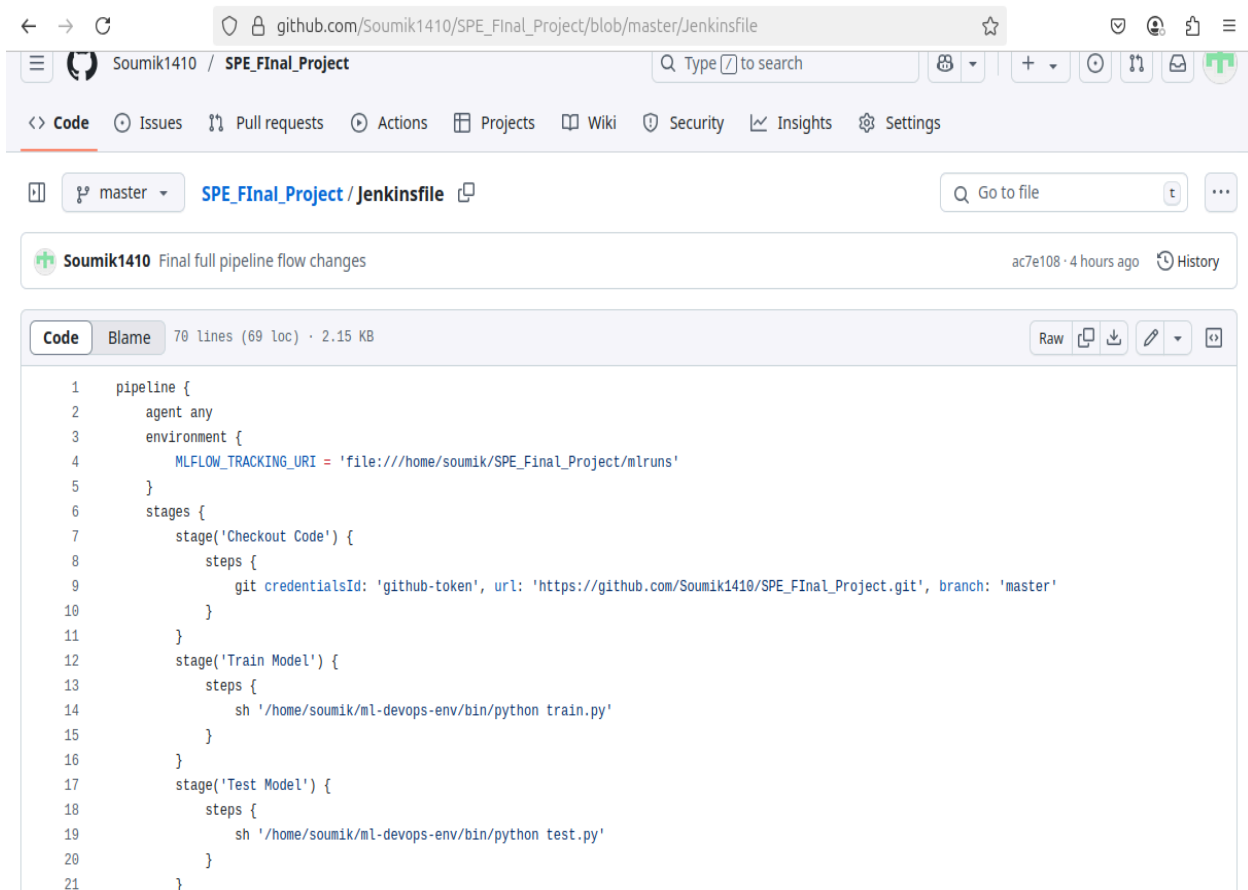
Script Path ?

Jenkinsfile

☒ Lightweight checkout ?

This creates a new Jenkins pipeline job that is triggered by webhooks sent from GitHub. It specifies that the source of the pipeline script resides in a Git repository, with details such as the repository URL, credentials for access, the branch where the pipeline script is located, and the filename of

the script— in our case, Jenkinsfile. The Jenkinsfile itself defines the pipeline using declarative syntax, outlining each stage and the steps required within those stages to execute the tasks mentioned earlier.



```
1 pipeline {
2   agent any
3   environment {
4     MLFLOW_TRACKING_URI = 'file:///home/soumik/SPE_Final_Project/mlruns'
5   }
6   stages {
7     stage('Checkout Code') {
8       steps {
9         git credentialsId: 'github-token', url: 'https://github.com/Soumik1410/SPE_Final_Project.git', branch: 'master'
10      }
11    }
12    stage('Train Model') {
13      steps {
14        sh '/home/soumik/ml-devops-env/bin/python train.py'
15      }
16    }
17    stage('Test Model') {
18      steps {
19        sh '/home/soumik/ml-devops-env/bin/python test.py'
20      }
21    }
22  }
```

For the first step i.e. checking out the latest code from a Github repository, the credentials needed were added to Jenkins credential manager. A Personal Access Token was created in Github, and that was stored in Jenkins as a username with password passing the PAT in place of the user's actual password.

Global credentials (unrestricted) [+ Add Credentials](#)

Credentials that should be available irrespective of domain specification to requirements matching.

ID	Name	Kind	Description
dockerhub-creds	soumik1410/*****	Username with password	
github-token	Soumik1410/*****	Username with password	

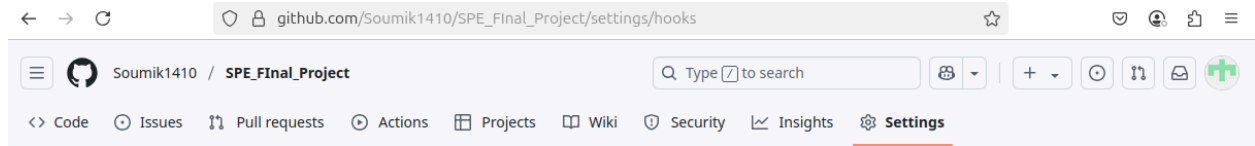
Icon: S M L

For the training and test stages, we needed all the libraries that were installed in the virtual environment used for development and as such first that environment is activated by using that specific environment's python interpreter instead of the system wise python interpreter located at `/usr/bin/python`.

MLFlow also normally stores all metrics and parameters and artifacts in a local folder 'mlruns' relative to the execution context of the train and test scripts. However, for these runs to show in the UI, they need to store them in a single location specified by the `MLFLOW_TRACKING_URI` environment variable.

Thus, the same is set in jenkins to ensure that different mlruns folders are not used to store the artifacts for different builds of the job, but instead they are all stored in the single specified location.

Now to trigger the jenkins pipeline whenever a push was made to the Github repository, a webhook was set up in Github that is delivered to jenkins everytime a push is made.



General

Access

Collaborators

Moderation options

Code and automation

Branches

Tags

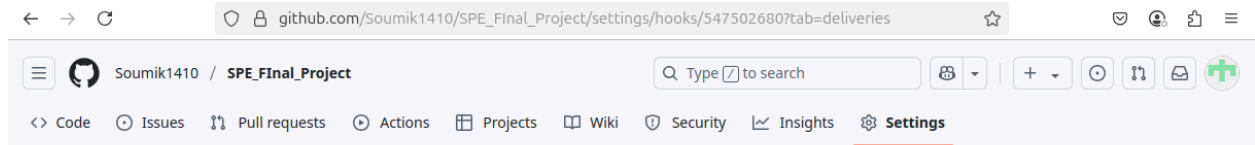
Webhooks

Webhooks allow external services to be notified when certain events happen. When the specified events happen, we'll send a POST request to each of the URLs you provide. Learn more in our [Webhooks Guide](#).

✓ [https://a7d7-49-207-60-154.ngrok... \(push\)](https://a7d7-49-207-60-154.ngrok...)

EditDelete

Last delivery was successful.



General

Access

Collaborators

Moderation options

Code and automation

Branches

Tags

Rules

Actions

Models

Webhooks

Environments

Codespaces

Pages

Webhooks / Manage webhook

SettingsRecent Deliveries

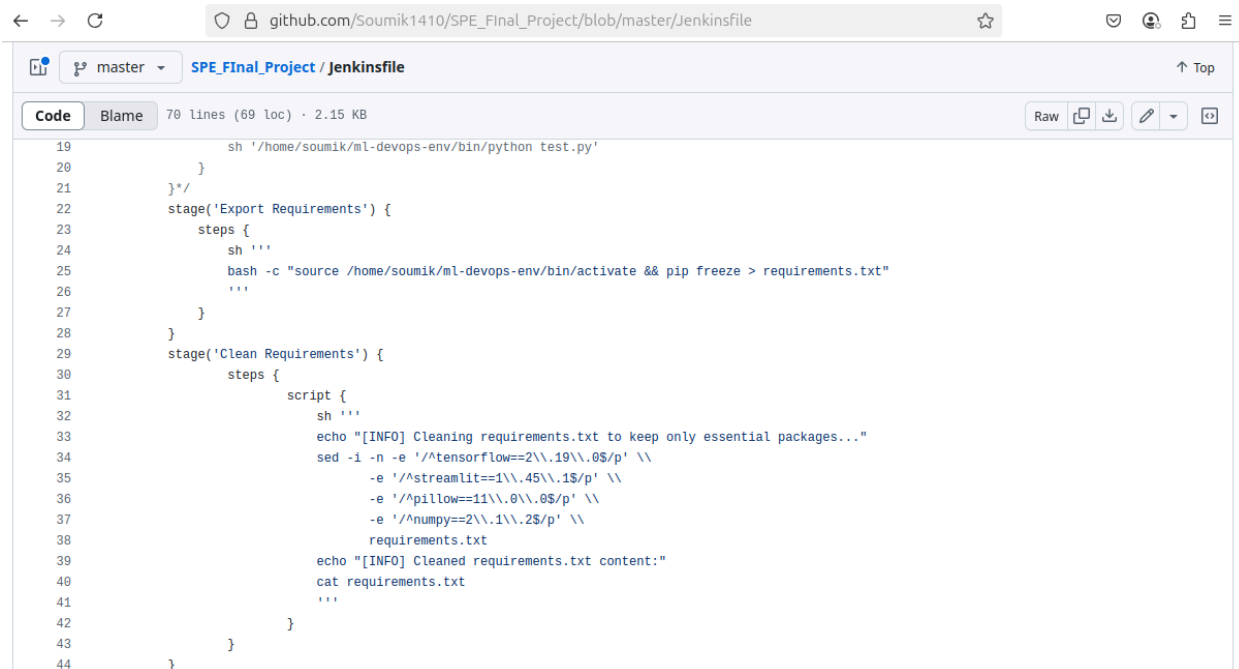
✓	7c5fa128-3641-11f0-958a-18a090d96c2a	push	2025-05-21 18:15:34	...
✓	62f6bc26-3641-11f0-9296-d1f4c411eddf	push	2025-05-21 18:14:52	...
✓	f25a4c36-3617-11f0-8ec9-8a7299c44bca	push	2025-05-21 13:18:13	...
✓	63bfccba-3614-11f0-9e36-e5bd97b8427a	push	2025-05-21 12:52:46	...
✓	194ffa2a-3613-11f0-9c97-5094f53f0426	push	2025-05-21 12:43:31	...
✓	097b626a-360e-11f0-8fe9-d8125ad2f405	push	2025-05-21 12:07:17	...
✓	03dc5ab4-360c-11f0-8256-b9506a203c9c	push	2025-05-21 11:52:49	...
✓	49276836-360a-11f0-896b-f2604d1ea0dd	push	2025-05-21 11:40:26	...

To enable the delivery of the webhook, the Jenkins instance running on my localhost needed to be exposed to the internet through a secure HTTP tunnel. This was achieved by using ngrok, which exposes port 8080 and generates a public URL. This `https://<ngrok URL>/github-webhook`, serves as the endpoint where the webhooks are delivered.

4. Dockerizing the Application :

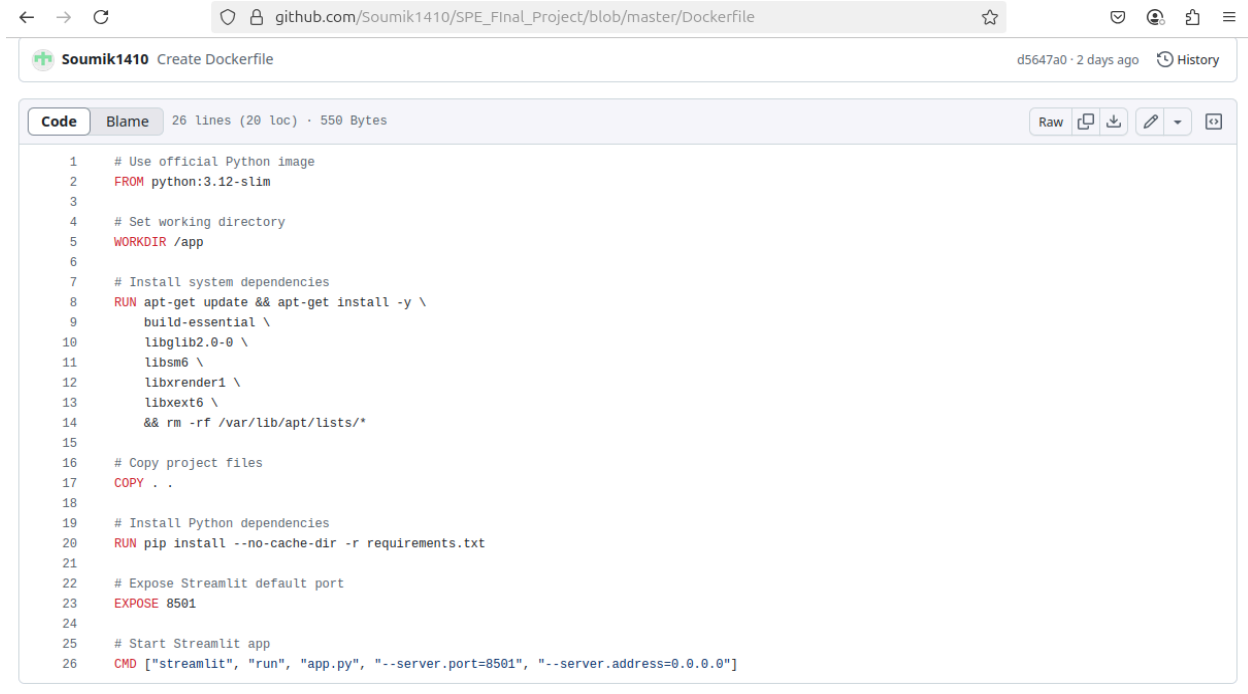
The next step in the project is to build an image that includes my inference script and the Streamlit-based frontend (defined in `app.py`). This image will be built on top of a base Python image, incorporating all necessary dependencies, the trained model, and any other saved artifacts required for deployment.

The next stage added to the Jenkins pipeline focuses on generating a `requirements.txt` file that lists all the environment's dependencies. However, since the environment is quite bloated due to the presence of training and testing scripts in the same repository, we need to prune the dependencies to include only the essential ones required to run `app.py` and `infer.py`. To achieve this, we added an additional stage in Jenkins to filter out all unnecessary dependencies from `requirements.txt`, leaving only Streamlit, TensorFlow, NumPy, Pandas, and their transitive dependencies. This helps reduce the overall size of the image.



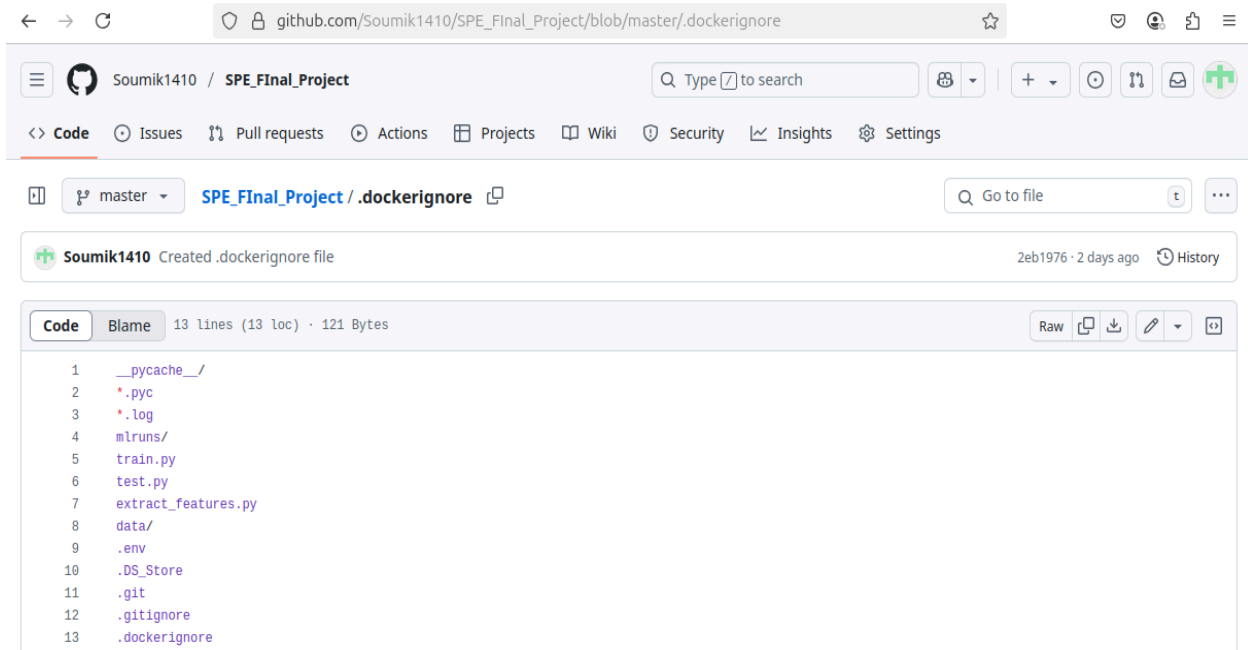
```
19         sh '/home/soumik/ml-devops-env/bin/python test.py'
20     }
21 }*/
22 stage('Export Requirements') {
23     steps {
24         sh '''
25         bash -c "source /home/soumik/ml-devops-env/bin/activate && pip freeze > requirements.txt"
26         '''
27     }
28 }
29 stage('Clean Requirements') {
30     steps {
31         script {
32             sh '''
33             echo "[INFO] Cleaning requirements.txt to keep only essential packages..."
34             sed -i -n -e '/^tensorflow==2\\.19\\.0$/p' \\
35                 -e '/^streamlit==1\\.45\\.1$/p' \\
36                 -e '/^pillow==11\\.0\\.0$/p' \\
37                 -e '/^numpy==2\\.1\\.2$/p' \\
38                 requirements.txt
39             echo "[INFO] Cleaned requirements.txt content:"
40             cat requirements.txt
41             '''
42         }
43     }
44 }
```

Then finally to create our docker image, we design our Dockerfile and .dockerignore files that Docker will use to build the image and also a new stage added to the Jenkins pipeline to trigger the docker build command.



The screenshot shows a GitHub web interface for a repository named 'Soumik1410' with the file 'Dockerfile' selected. The browser address bar shows 'github.com/Soumik1410/SPE_Final_Project/blob/master/Dockerfile'. The file content is as follows:

```
1 # Use official Python image
2 FROM python:3.12-slim
3
4 # Set working directory
5 WORKDIR /app
6
7 # Install system dependencies
8 RUN apt-get update && apt-get install -y \
9     build-essential \
10     libgl2.0-0 \
11     libsm6 \
12     libxrender1 \
13     libxext6 \
14     && rm -rf /var/lib/apt/lists/*
15
16 # Copy project files
17 COPY . .
18
19 # Install Python dependencies
20 RUN pip install --no-cache-dir -r requirements.txt
21
22 # Expose Streamlit default port
23 EXPOSE 8501
24
25 # Start Streamlit app
26 CMD ["streamlit", "run", "app.py", "--server.port=8501", "--server.address=0.0.0.0"]
```



The screenshot shows a GitHub web interface for the same repository 'Soumik1410' with the file '.dockerignore' selected. The browser address bar shows 'github.com/Soumik1410/SPE_Final_Project/blob/master/.dockerignore'. The file content is as follows:

```
1 __pycache__/
2 *.pyc
3 *.log
4 mlruns/
5 train.py
6 test.py
7 extract_features.py
8 data/
9 .env
10 .DS_Store
11 .git
12 .gitignore
13 .dockerignore
```



```
40         cat requirements.txt
41         '''
42     }
43 }
44 }
45 stage("Build Docker Image") {
46     steps {
47         script {
48             sh 'docker build -t imagecaptioner_mt2024153:latest .'
49             echo 'Docker Image successfully created.'
50         }
51     }
52 }
```

Docker will build the image according to the instructions given in the Dockerfile i.e. using python:3.12-slim as the base image, installing system dependencies, copying everything from jenkins workspace current directory to docker build (for the saved model and artifacts), then installing required python dependencies, exposing streamlit's default port, then finally running the streamlit point as entry command.

.dockerfile lists the files and folders to not be copied when executing the COPY . . command in Dockerfile

Once the build is complete, it is pushed to the Docker Hub repository. To facilitate this, a new stage is added to the Jenkins pipeline, and Docker Hub credentials are securely stored in Jenkins' credential manager for seamless authentication during the push process.

```
53     stage('Pushing Docker Image to Hub') {
54     steps {
55         script {
56             withDockerRegistry([credentialsId: 'dockerhub-creds', url: 'https://index.docker.io/v1/']) {
57                 sh 'docker tag imagecaptioner_mt2024153:latest soumik1410/imagecaptioner_mt2024153:latest'
58                 sh 'docker push soumik1410/imagecaptioner_mt2024153:latest'
59             }
60         }
61     }
62 }
```

←

→

↺


a7d7-49-207-60-154.ngrok-free.app/manage/credentials/store/system/domain/_/

☆

🔒

👤

📄

Jenkins

🔍

🔔

🔒

🔴

admin

▼

🔑



log out

Dashboard > Manage Jenkins > Credentials > System > Global credentials (unrestricted) >

Global credentials (unrestricted)

[+ Add Credentials](#)

Credentials that should be available irrespective of domain specification to requirements matching.

ID	Name	Kind	Description
 dockerhub-creds	soumik1410/*****	Username with password	

←

→

↺

hub.docker.com/repositories/soumik1410

☆

🔒

👤


📄

☰

New

Introducing Docker Hardened Images - Learn More →

×

dockerhub

Explore

My Hub

🔍

Search Docker Hub

CtrlK


🔔

🔒

🔄

☰

S

soumik1410
Docker Personal

▼

🔧

Repositories

👤

Collaborations

⚙️

Settings

▼

Default privacy

Notifications

💰

Billing

📄

Usage

▼

Pulls

Storage

Repositories

All repositories within the soumik1410 namespace.

🔍

Search by repository name

All content

▼

Create a repository

Name	Last Pushed	Contains	Visibility	Scout
soumik1410/imagecaptioner_mt2024153	about 18 hours ago	IMAGE	Public	Inactive
soumik1410/calculator_mt2024153	about 2 months ago	IMAGE	Public	Inactive

1-2 of 2

◀ ▶

The screenshot shows the Docker Hub interface for the repository `soumik1410/imagecaptioner_mt2024153`. The left sidebar contains navigation links for Repositories, Collaborations, Settings, Default privacy, Notifications, Billing, Usage, Pulls, and Storage. The main content area shows the repository details, including the name, last pushed time, and repository size. The 'Tags' section displays a table of tags, and the 'Docker commands' section provides a public view button and a command to push a new tag.

Tag	OS	Type	Pulled	Pushed
latest	linux	Image	less than 1 day	about 19 hours

5. Ansible - Deploying the application to a Kubernetes cluster (with Ansible Roles for modularity) :

Once the build has been successfully pushed to Docker Hub, the next step is to deploy it to the target environment. The deployment will be carried out on a Kubernetes cluster running locally using Minikube. The deployment tasks involve several key actions: pulling the latest build from Docker Hub, gracefully stopping the existing Minikube instance, starting up Minikube, setting the kubectl context to Minikube, creating the necessary namespace, and applying the Kubernetes manifests for both the application deployment and service.

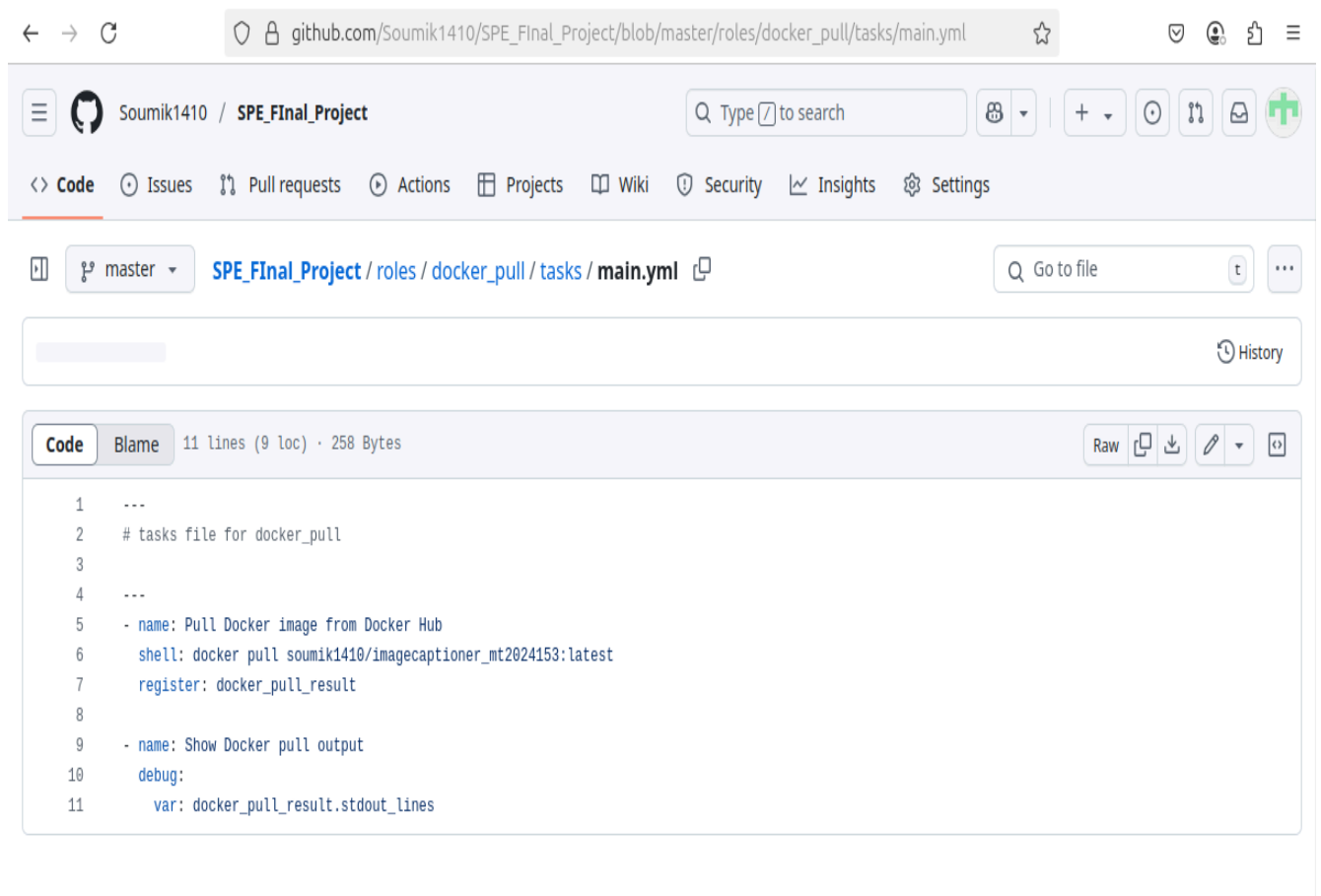
To automate these steps, we utilize an Ansible playbook that defines each task using the appropriate Ansible modules and specifies the desired states or actions. This ensures a smooth, repeatable deployment process on the local Minikube cluster.

We utilize Ansible roles to break up all the deployment tasks into 2 modules, 1 for pulling the latest docker image from Docker Hub and another for handling deployment of the application to the Kubernetes cluster.

The 2 roles created were `docker_pull` and `k8s_deploy`.

Folder structure for this included a roles folder with these 2 subdirectories and both followed a common structure with subfolders for tasks, handlers, tests etc. We have only used tasks in our project, modifying the `main.yml` files for the 2 roles to perform the tasks mentioned earlier.

roles/docker_pull/tasks/main.yml :



The screenshot shows a web browser displaying a GitHub repository page. The address bar shows the URL: `github.com/Soumik1410/SPE_Final_Project/blob/master/roles/docker_pull/tasks/main.yml`. The repository name is `SPE_Final_Project` by user `Soumik1410`. The file path is `roles / docker_pull / tasks / main.yml`. The file content is displayed in a code editor with line numbers 1 through 11. The content is an Ansible task file for the `docker_pull` role.

```
1  ---
2  # tasks file for docker_pull
3
4  ---
5  - name: Pull Docker image from Docker Hub
6    shell: docker pull soumik1410/imagecaptioner_mt2024153:latest
7    register: docker_pull_result
8
9  - name: Show Docker pull output
10    debug:
11      var: docker_pull_result.stdout_lines
```

roles/k8s_deploy/tasks/main.yaml :

← → ↻ github.com/Soumik1410/SPE_Final_Project/blob/master/roles/k8s_deploy/tasks/main.yaml ☆

☰ Soumik1410 / SPE_Final_Project 🔍 Type to search 🗄️ + 🔍 🔗 📧 +

<> Code 🔍 Issues 🔗 Pull requests 🔄 Actions 📁 Projects 📖 Wiki 🛡️ Security 🔍 Insights ⚙️ Settings

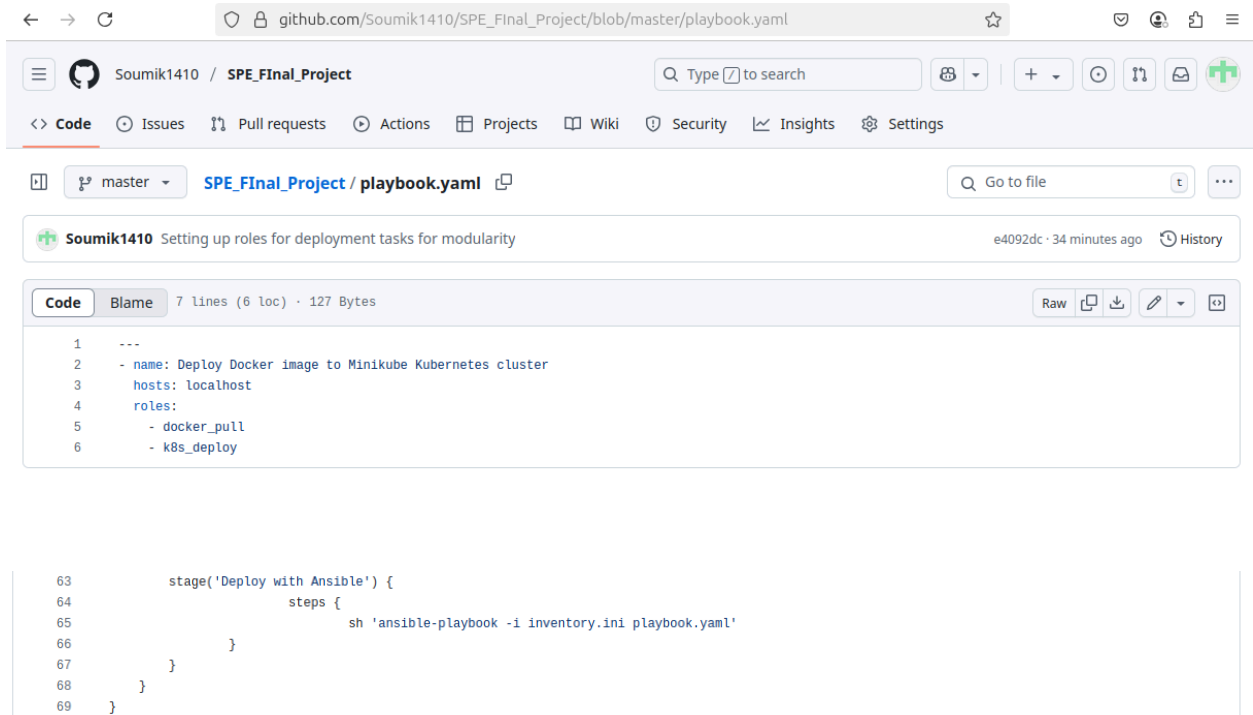
📄 master SPE_Final_Project / roles / k8s_deploy / tasks / main.yaml 🔍 Go to file 🔍 ...

🕒 History

Code Blame 40 lines (32 loc) · 1.05 KB Raw 📄 ⬇️ ✎ ...

```
1  |--
2  # tasks file for roles/k8s_deploy
3
4  - name: Stop existing Minikube cluster if running
5    shell: minikube stop
6    ignore_errors: yes
7
8  - name: Delete existing Minikube cluster if exists
9    shell: minikube delete
10   ignore_errors: yes
11
12 - name: Start Minikube if not already running
13   shell: |
14     minikube status || minikube start --driver=docker
15   register: minikube_status
16   changed_when: "'Running' not in minikube_status.stdout"
17
18 - name: Set kubectl context to Minikube
19   shell: minikube update-context
20
21 - name: Create namespace if not present
22   kubernetes.core.k8s:
23     api_version: v1
24     kind: Namespace
25     name: image-caption-app-mt2024153
26     state: present
27
28 - name: Apply Kubernetes deployment
29   kubernetes.core.k8s:
30     kubeconfig: ~/.kube/config
31     state: present
32     definition: "{{ lookup('file', 'deployment.yaml') }}"
33
34 - name: Apply Kubernetes service (if any)
35   kubernetes.core.k8s:
36     kubeconfig: ~/.kube/config
37     state: present
38     definition: "{{ lookup('file', 'service.yaml') }}"
39     when: lookup('file', 'service.yaml', errors='ignore') is not none
```

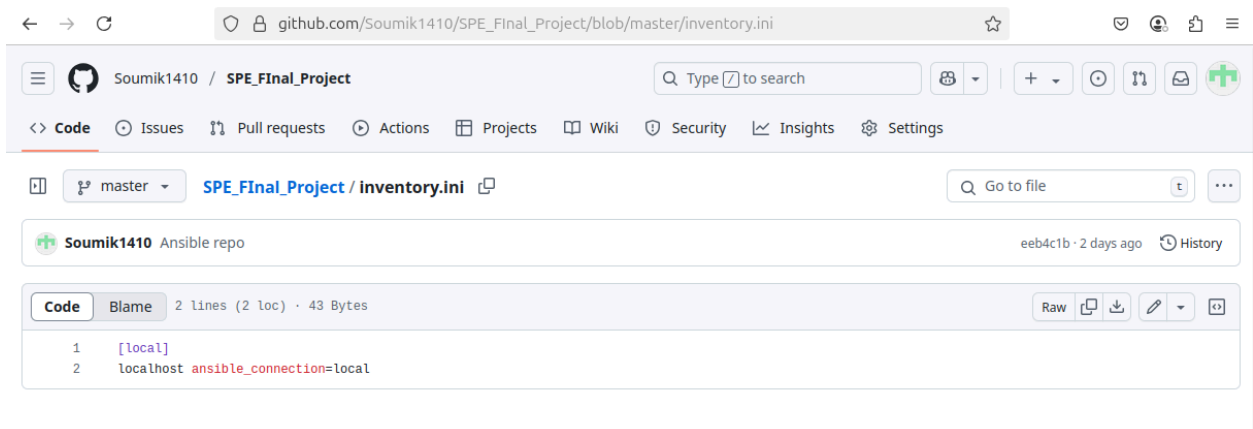

The main playbook file simply lists these roles and a new stage is added to the Jenkins pipeline that executes this main playbook file



```
---
- name: Deploy Docker image to Minikube Kubernetes cluster
  hosts: localhost
  roles:
    - docker_pull
    - k8s_deploy
```

```
63     stage('Deploy with Ansible') {
64         steps {
65             sh 'ansible-playbook -i inventory.ini playbook.yaml'
66         }
67     }
68 }
69 }
```

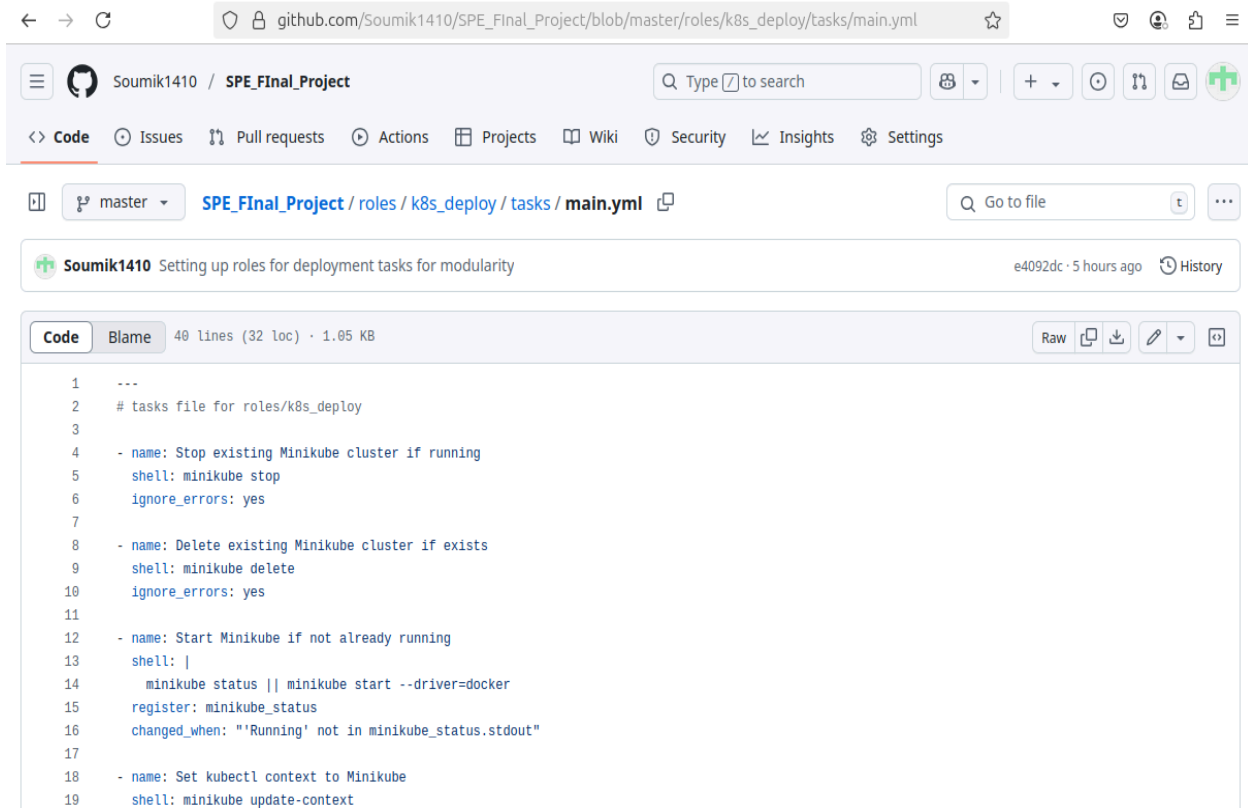
Ansible inventory file used for this deployment is simple as deployment target is localhost



```
[local]
localhost ansible_connection=local
```

Our deployment target is a kubernetes cluster deployed on my localhost node using minikube. The Ansible role k8s_deploy is used to set up the Kubernetes cluster and deploy our image captioning application.

First, it gracefully shuts down any existing minikube instance that might be running and starts up a new minikube instance, setting kubectl context to the kubernetes cluster deployed by the fresh minikube instance so that the cluster can be managed easily with kubectl commands without needing to specify cluster everytime.



The screenshot shows a GitHub repository page for 'Soumik1410 / SPE_Final_Project'. The file path is 'roles/k8s_deploy/tasks/main.yml'. The commit message is 'Setting up roles for deployment tasks for modularity'. The file is 40 lines (32 loc) and 1.05 KB. The code content is as follows:

```
1 ---
2 # tasks file for roles/k8s_deploy
3
4 - name: Stop existing Minikube cluster if running
5   shell: minikube stop
6   ignore_errors: yes
7
8 - name: Delete existing Minikube cluster if exists
9   shell: minikube delete
10  ignore_errors: yes
11
12 - name: Start Minikube if not already running
13   shell: |
14     minikube status || minikube start --driver=docker
15   register: minikube_status
16   changed_when: "'Running' not in minikube_status.stdout"
17
18 - name: Set kubectl context to Minikube
19   shell: minikube update-context
```

Then, we use kubectl commands to create a custom namespace for our application called image-caption-app-mt2024153 and apply our deployment and service Kubernetes manifests which are used to specify the number of replicas, the build deployed and network access to our application.

Ansible module used for these is `kubernetes.core.k8s` which can be installed with `ansible-galaxy collection install kubernetes.core`

First, a separate namespace is created for our application called image-caption-app-mt2024153.

```
21 - name: Create namespace if not present
22   kubernetes.core.k8s:
23     api_version: v1
24     kind: Namespace
25     name: image-caption-app-mt2024153
26     state: present
```

Then, we apply a deployment.yaml and service.yaml in that namespace that deploys our application's docker image in a container inside a Kubernetes pod and sets up the network access to it. State:present with kubernetes.core.k8s module is equivalent to kubectl apply.

```
28 - name: Apply Kubernetes deployment
29   kubernetes.core.k8s:
30     kubeconfig: ~/.kube/config
31     state: present
32     definition: "{{ lookup('file', 'deployment.yaml') }}"
33
34 - name: Apply Kubernetes service (if any)
35   kubernetes.core.k8s:
36     kubeconfig: ~/.kube/config
37     state: present
38     definition: "{{ lookup('file', 'service.yaml') }}"
39     when: lookup('file', 'service.yaml', errors='ignore') is not none
```

deployment.yaml :

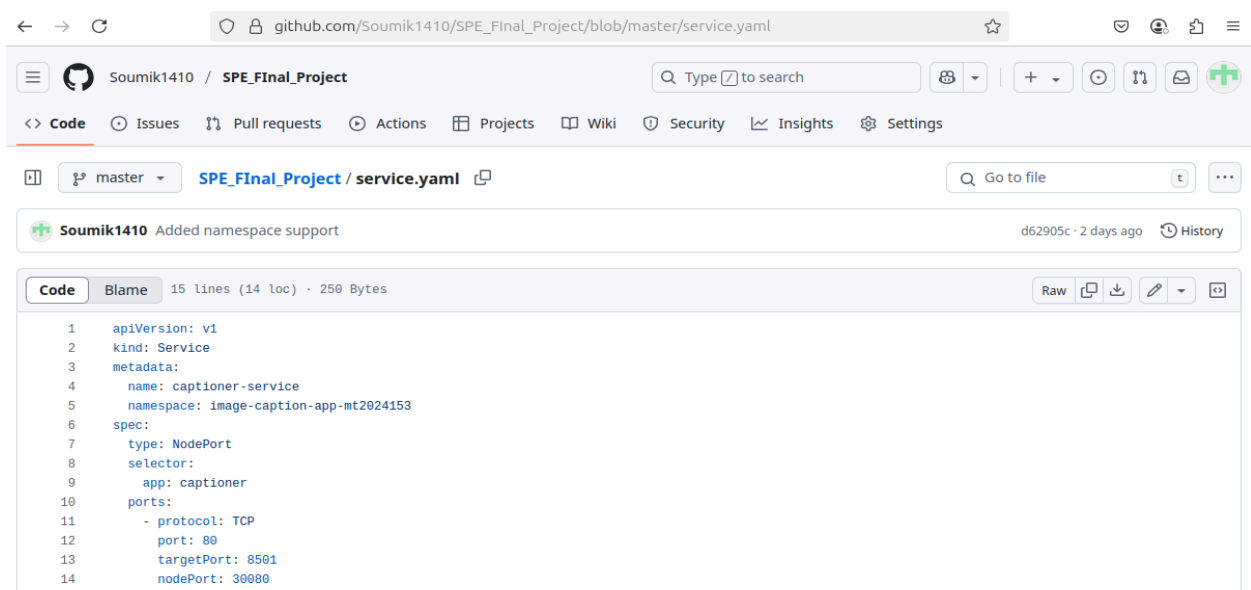
The screenshot shows a GitHub repository page for the file `deployment.yaml` in the repository `SPE_Final_Project`. The file is 21 lines long (20 loc) and 403 bytes. The content of the file is a Kubernetes Deployment manifest:

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: image-caption-app
5    namespace: image-caption-app-mt2024153
6  spec:
7    replicas: 1
8    selector:
9      matchLabels:
10       app: captioner
11    template:
12      metadata:
13       labels:
14         app: captioner
15      spec:
16       containers:
17         - name: captioner
18           image: soumik1410/imagecaptioner_mt2024153:latest
19           ports:
20             - containerPort: 8501
```

This configuration defines a Kubernetes Deployment for the Image Captioning App, called image-caption-app, intended to run within the image-caption-app-mt2024153 namespace. It specifies a single replica (i.e., one Pod) for the application. Both the application and the Pod are labeled with captioner to facilitate identification, enabling users to easily determine which application is running on which Pod, as well as to manage the associated Pods within the Deployment.

The YAML further details that the Deployment will launch a single container within the Pod, which will run the Docker image for the Image Captioning App. This container will expose port 8501, the required port for accessing the Streamlit application hosted on the Pod.

service.yaml



```
1  apiVersion: v1
2  kind: Service
3  metadata:
4    name: captioner-service
5    namespace: image-caption-app-mt2024153
6  spec:
7    type: NodePort
8    selector:
9      app: captioner
10   ports:
11     - protocol: TCP
12       port: 80
13       targetPort: 8501
14       nodePort: 30080
```

This YAML specifies a Kubernetes Service called captioner-service, intended to run within the namespace image-caption-app-mt2024153. It is a NodePort service, which is used to allow external HTTP connections to containers deployed within this namespace. This YAML defines a Kubernetes Service named captioner-service, which is deployed within the image-caption-app-mt2024153 namespace. The service is of type

NodePort, allowing external HTTP traffic to access containers running in the Kubernetes cluster. It uses the app: captioner label selector to identify and route traffic to the appropriate Pods.

The service maps node port 30080 to the container's port 8501 — the default port used by Streamlit. This setup allows external clients to connect to the Streamlit application running inside the Pods.

Once the deployment is complete we can verify the status using kubectl commands.

```
jenkins@soumik-VirtualBox:~$ minikube status
minikube
type: Control Plane
host: Running
kubelet: Running
apiserver: Running
kubeconfig: Configured

jenkins@soumik-VirtualBox:~$
```

```
jenkins@soumik-VirtualBox:~$ kubectl get ns
NAME                STATUS    AGE
default             Active   22m
image-caption-app-mt2024153  Active   22m
kube-node-lease     Active   22m
kube-public         Active   22m
kube-system         Active   22m
jenkins@soumik-VirtualBox:~$
```

```
jenkins@soumik-VirtualBox:~$ kubectl get pods -n image-caption-app-mt2024153
NAME                                READY   STATUS    RESTARTS   AGE
fluent-bit-c4zh6                   1/1     Running   0           17m
grafana-8559bcb694-2xpd           1/1     Running   0           17m
image-caption-app-756c979f57-zmjnh 1/1     Running   0           18m
loki-87d578f98-c22dr              1/1     Running   0           17m
jenkins@soumik-VirtualBox:~$
```

```
jenkins@soumik-VirtualBox:~$ kubectl describe pod image-caption-app-756c979f57-zmjnh -n image-caption-app-mt2024153
Name:          image-caption-app-756c979f57-zmjnh
Namespace:     image-caption-app-mt2024153
Priority:       0
Service Account: default
Node:          minikube/192.168.58.2
Start Time:    Thu, 22 May 2025 15:42:58 +0530
Labels:        app=captioner
               pod-template-hash=756c979f57
Annotations:   <none>
Status:        Running
IP:            10.244.0.5
IPs:
  IP:          10.244.0.5
Controlled By: ReplicaSet/image-caption-app-756c979f57
Containers:
  captioner:
    Container ID:  docker://4019ff149aada3b8bed5b265fd142a395df318b46541688933f7e1dbd75f1169
    Image:          soumik1410/imagecaptioner_mt2024153:latest
    Image ID:       docker-pullable://soumik1410/imagecaptioner_mt2024153@sha256:f24895847a4e7558a26f4d598765c7b9f38460b39b3a9b86261e72e32c6bb19b
    Port:          8501/TCP
    Host Port:     0/TCP
    State:         Running
      Started:     Thu, 22 May 2025 15:43:01 +0530
    Ready:         True
    Restart Count: 0
    Limits:
      cpu: 500m
    Requests:
      cpu: 200m
```

```
Environment: <none>
Mounts:
  /var/run/secrets/kubernetes.io/serviceaccount from kube-api-access-2j4qd (ro)
Conditions:
  Type              Status
  PodReadyToStartContainers  True
  Initialized        True
  Ready              True
  ContainersReady    True
  PodScheduled       True
Volumes:
  kube-api-access-2j4qd:
    Type:              Projected (a volume that contains injected data from multiple sources)
    TokenExpirationSeconds: 3607
    ConfigMapName:      kube-root-ca.crt
    ConfigMapOptional:  <nil>
    DownwardAPI:        true
QoS Class:           Burstable
Node-Selectors:      <none>
Tolerations:         node.kubernetes.io/not-ready:NoExecute op=Exists for 300s
                     node.kubernetes.io/unreachable:NoExecute op=Exists for 300s
Events:
  Type    Reason      Age   From              Message
  ----    -
  Normal  Scheduled   19m   default-scheduler Successfully assigned image-caption-app-mt2024153/image-caption-app-756c979f57-zmjnh to minikube
  Normal  Pulling     19m   kubelet           Pulling image "soumik1410/imagecaptioner_mt2024153:latest"
  Normal  Pulled      19m   kubelet           Successfully pulled image "soumik1410/imagecaptioner_mt2024153:latest" in 2.137s (2.137s including waiting). Image size: 3039320995 bytes.
  Normal  Created     19m   kubelet           Created container: captioner
  Normal  Started     19m   kubelet           Started container captioner
jenkins@soumik-VirtualBox:~$
```

```
jenkins@soumik-VirtualBox:~$ kubectl get deployment -n image-caption-app-mt2024153
NAME          READY   UP-TO-DATE   AVAILABLE   AGE
grafana       1/1     1             1           21m
image-caption-app 1/1     1             1           26m
loki          1/1     1             1           21m
jenkins@soumik-VirtualBox:~$ kubectl get svc -n image-caption-app-mt2024153
NAME          TYPE        CLUSTER-IP   EXTERNAL-IP   PORT(S)          AGE
captioner-service NodePort    10.104.28.183 <none>        80:30080/TCP     26m
grafana       NodePort    10.107.126.227 <none>        3000:30030/TCP   21m
loki          ClusterIP   10.99.136.160 <none>        3100/TCP         21m
jenkins@soumik-VirtualBox:~$
```

As we can see, all our deployment activities went successful. Here, we can see a couple more pods and deployments, these were brought up to introduce Horizontal Pod Autoscaling (HPA) and logging and monitoring in the pipeline.

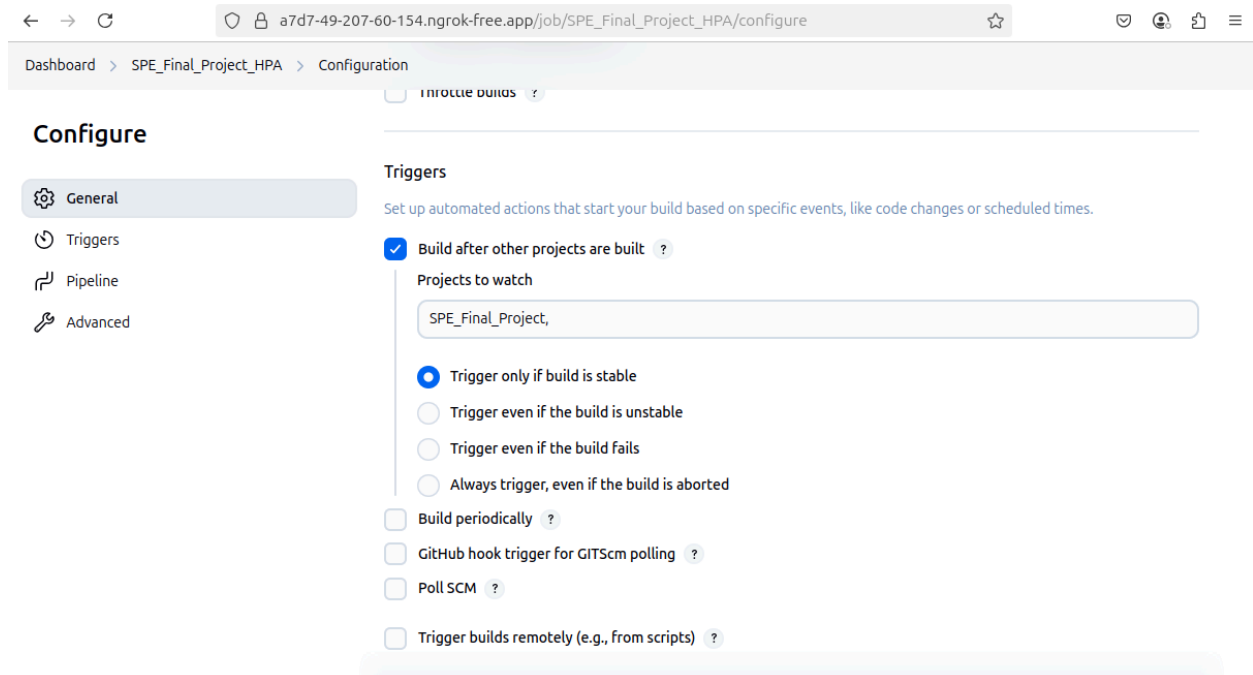
6. Horizontal Pod Autoscaler :

HPA (Horizontal Pod Autoscaler) is a Kubernetes resource that automatically scales the number of Pod replicas in a Deployment based on observed resource usage, such as CPU or memory, or on custom metrics.

For this to work in our cluster, we need to first bring up the metrics-server which measures metrics such as CPU usage or memory usage of pods. We need to set requests of CPU and memory usage for our image-caption-app Deployment. Then in our HPA configuration, we need to set the resource that we want to monitor, the target utilization (computed as a percentage of the requested value) and min and max replicas that we desire. Then based on current utilization vs target utilization, HPA will automatically scale up and scale down the number of replicas to ensure average resource utilization is close to the target threshold.

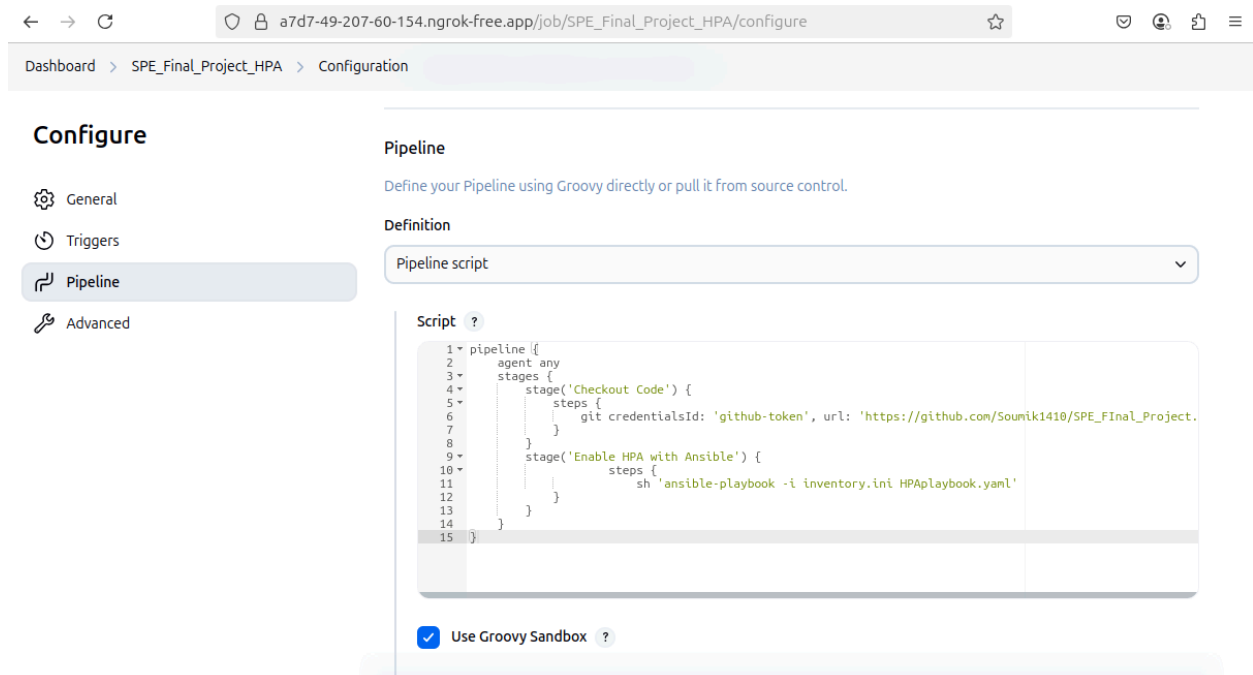
We define another Ansible playbook called HPAplaybook.yaml that lists all these tasks using the kubernetes.core.k8s module. This playbook is executed automatically once the application deployment is completed, by a Jenkins pipeline which checks out the latest version of this playbook from our Github repository and then executes this playbook. This is a separate

pipeline from our first Jenkins pipeline that handled automated model training, testing, containerization and deployment of the application. This was done to modularize and independently test the HPA without needing to train the model time and time again. This second pipeline is automatically triggered on successful completion of the first pipeline.



The screenshot shows the Jenkins Configuration page for the job 'SPE_Final_Project_HPA'. The 'Triggers' tab is selected in the left sidebar. The main content area is titled 'Triggers' and includes a description: 'Set up automated actions that start your build based on specific events, like code changes or scheduled times.' The configuration options are as follows:

- ☒ Build after other projects are built ?
 - Projects to watch: SPE_Final_Project,
 - ☒ Trigger only if build is stable
 - ☐ Trigger even if the build is unstable
 - ☐ Trigger even if the build fails
 - ☐ Always trigger, even if the build is aborted
- ☐ Build periodically ?
- ☐ GitHub hook trigger for GITScm polling ?
- ☐ Poll SCM ?
- ☐ Trigger builds remotely (e.g., from scripts) ?



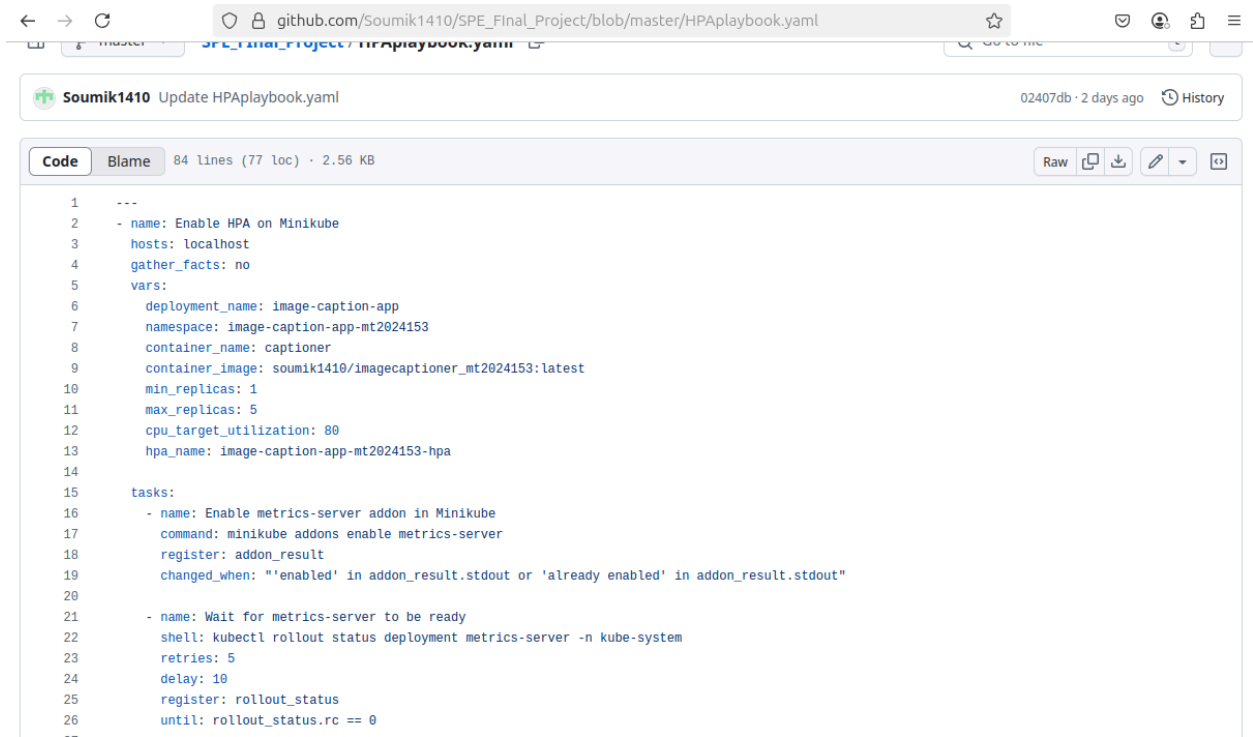
The screenshot shows the Jenkins Configuration page for the job 'SPE_Final_Project_HPA', with the 'Pipeline' tab selected in the left sidebar. The main content area is titled 'Pipeline' and includes a description: 'Define your Pipeline using Groovy directly or pull it from source control.' The configuration options are as follows:

- Definition: Pipeline script (selected from a dropdown menu)
- Script ?

```
1 pipeline {  
2   agent any  
3   stages {  
4     stage('Checkout Code') {  
5       steps {  
6         git credentialsId: 'github-token', url: 'https://github.com/Soumik1410/SPE_Final_Project.  
7       }  
8     }  
9     stage('Enable HPA with Ansible') {  
10      steps {  
11        sh 'ansible-playbook -i inventory.ini HPAPlaybook.yaml'  
12      }  
13    }  
14  }  
15 }
```
- ☒ Use Groovy Sandbox ?

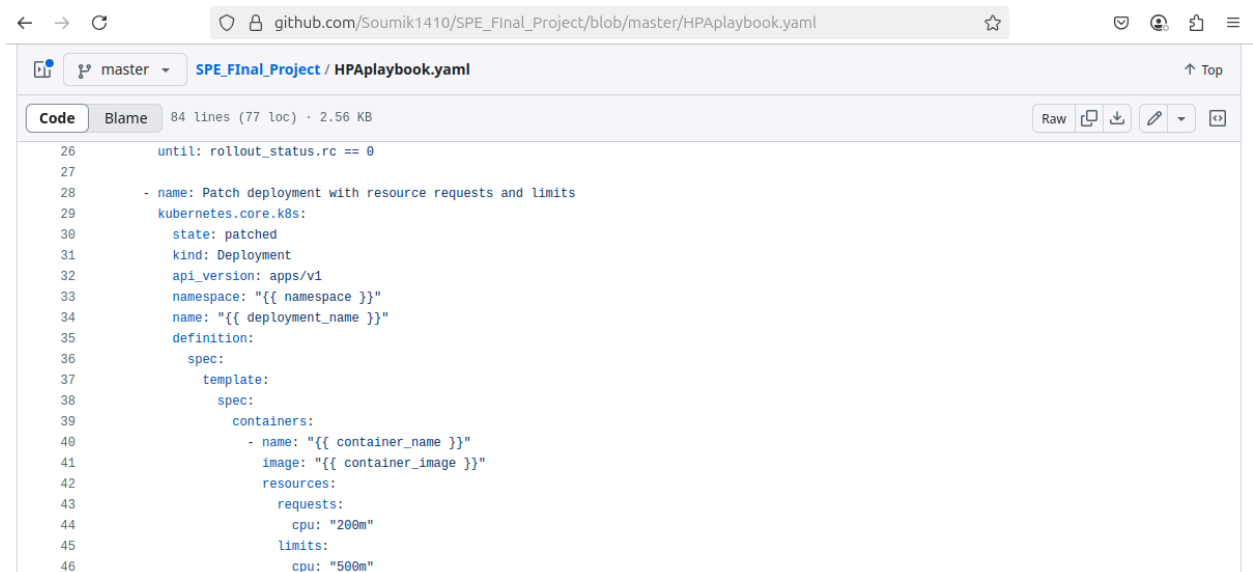
HPAplaybook.yaml tasks :

Starting metrics server and waiting for it to be ready



```
1 ---
2 - name: Enable HPA on Minikube
3   hosts: localhost
4   gather_facts: no
5   vars:
6     deployment_name: image-caption-app
7     namespace: image-caption-app-mt2024153
8     container_name: captioner
9     container_image: soumik1410/imagecaptioner_mt2024153:latest
10    min_replicas: 1
11    max_replicas: 5
12    cpu_target_utilization: 80
13    hpa_name: image-caption-app-mt2024153-hpa
14
15  tasks:
16    - name: Enable metrics-server add-on in Minikube
17      command: minikube addons enable metrics-server
18      register: addon_result
19      changed_when: "'enabled' in addon_result.stdout or 'already enabled' in addon_result.stdout"
20
21    - name: Wait for metrics-server to be ready
22      shell: kubectl rollout status deployment metrics-server -n kube-system
23      retries: 5
24      delay: 10
25      register: rollout_status
26      until: rollout_status.rc == 0
27
28  ~
```

Patch our application deployment to have resource requests and limits



```
26    until: rollout_status.rc == 0
27
28  - name: Patch deployment with resource requests and limits
29    kubernetes.core.k8s:
30      state: patched
31      kind: Deployment
32      api_version: apps/v1
33      namespace: "{{ namespace }}"
34      name: "{{ deployment_name }}"
35      definition:
36        spec:
37          template:
38            spec:
39              containers:
40                - name: "{{ container_name }}"
41                  image: "{{ container_image }}"
42                  resources:
43                    requests:
44                      cpu: "200m"
45                    limits:
46                      cpu: "500m"
```

Configure HPA

```
← → ↺ github.com/Soumik1410/SPE_Final_Project/blob/master/HPAplaybook.yaml ☆ 🔒 👤 📄 ☰
SPE_Final_Project / HPAplaybook.yaml ↑ Top
Code Blame 84 lines (77 loc) · 2.56 KB Raw 📄 ⬇️ ✎ ⌵ 🗨️
47
48 - name: Create or update HorizontalPodAutoscaler
49   kubernetes.core.k8s:
50     state: present
51     definition:
52       apiVersion: autoscaling/v2
53       kind: HorizontalPodAutoscaler
54       metadata:
55         name: "{{ hpa_name }}"
56         namespace: "{{ namespace }}"
57       spec:
58         scaleTargetRef:
59           apiVersion: apps/v1
60           kind: Deployment
61           name: "{{ deployment_name }}"
62         minReplicas: "{{ min_replicas }}"
63         maxReplicas: "{{ max_replicas }}"
64         metrics:
65         - type: Resource
66           resource:
67             name: cpu
68           target:
69             type: Utilization
70             averageUtilization: "{{ cpu_target_utilization }}"
```

Wait for HPA to start up and ready to start monitoring

```
72 - name: Wait for HPA to be created
73   shell: kubectl get hpa "{{ hpa_name }}" -n "{{ namespace }}"
74   retries: 5
75   delay: 5
76   register: hpa_status
77   until: hpa_status.rc == 0
78
79 - name: Show HPA status
80   command: kubectl describe hpa "{{ hpa_name }}" -n "{{ namespace }}"
81   register: hpa_describe
82
83 - debug:
84   msg: "{{ hpa_describe.stdout }}"
```

Once this is executed, we can verify the status using kubectl commands.

```
jenkins@soumik-VirtualBox:~$ kubectl rollout status deployment metrics-server -n kube-system
deployment "metrics-server" successfully rolled out
jenkins@soumik-VirtualBox:~$ kubectl get pods -n kube-system
NAME                                READY   STATUS    RESTARTS   AGE
coredns-668d6bf9bc-zjgtz            1/1     Running   0           4h35m
etcd-minikube                        1/1     Running   0           4h36m
kube-apiserver-minikube              1/1     Running   0           4h35m
kube-controller-manager-minikube     1/1     Running   0           4h36m
kube-proxy-kb4st                     1/1     Running   0           4h35m
kube-scheduler-minikube              1/1     Running   0           4h36m
metrics-server-7fbb699795-5zxvg      1/1     Running   0           4h35m
storage-provisioner                  1/1     Running   1 (4h35m ago) 4h35m
jenkins@soumik-VirtualBox:~$ kubectl top pods -n image-caption-app-mt2024153
NAME                                CPU(cores)   MEMORY(bytes)
fluent-bit-c4zh6                     98m          12Mi
grafana-8559bcb694-2xpcd             2m           41Mi
image-caption-app-756c979f57-zmjnh    0m           7Mi
loki-87d578f98-c22dr                 33m          280Mi
jenkins@soumik-VirtualBox:~$
```

```
jenkins@soumik-VirtualBox:~$ kubectl describe deployment image-caption-app -n image-caption-app-mt2024153
Name: image-caption-app
Namespace: image-caption-app-mt2024153
CreationTimestamp: Thu, 22 May 2025 15:38:24 +0530
Labels: <none>
Annotations: deployment.kubernetes.io/revision: 2
Selector: app=captioner
Replicas: 1 desired | 1 updated | 1 total | 1 available | 0 unavailable
StrategyType: RollingUpdate
MinReadySeconds: 0
RollingUpdateStrategy: 25% max unavailable, 25% max surge
Pod Template:
  Labels: app=captioner
  Containers:
    captioner:
      Image: soumik1410/imagecaptioner_mt2024153:latest
      Port: 8501/TCP
      Host Port: 0/TCP
      Limits:
        cpu: 500m
      Requests:
        cpu: 200m
      Environment: <none>
      Mounts: <none>
      Volumes: <none>
      Node-Selectors: <none>
      Tolerations: <none>
  Conditions:
    Type           Status  Reason
    ----           -
    Available       True    MinimumReplicasAvailable
    Progressing     True    NewReplicaSetAvailable
```

```
jenkins@soumik-VirtualBox:~$ kubectl get hpa -n image-caption-app-mt2024153
```

NAME	REFERENCE	TARGETS	MINPODS	MAXPODS	REPLICAS	AGE
image-caption-app-mt2024153-hpa	Deployment/image-caption-app	cpu: 0%/80%	1	5	1	4h36m

```
jenkins@soumik-VirtualBox:~$
```

7. Logging and Monitoring :