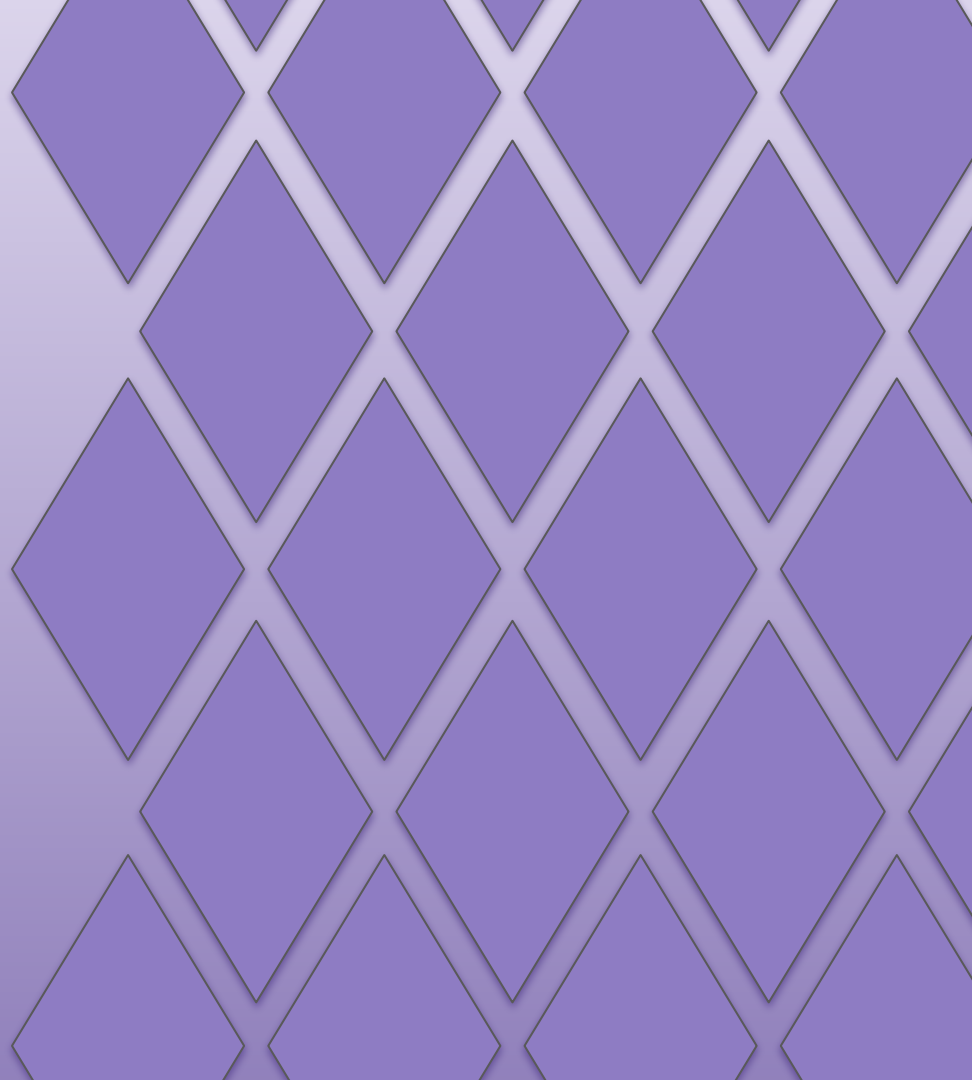


Project:
**Classification of Rakuten
e-commerce products**





Team members:

Karim Osman
&
Josipa Rupčić

Project overview:

- This project involves deploying three machine learning models for an e-commerce platform
- Goal of the project is to classify product descriptions and images into predefined categories

Model used:

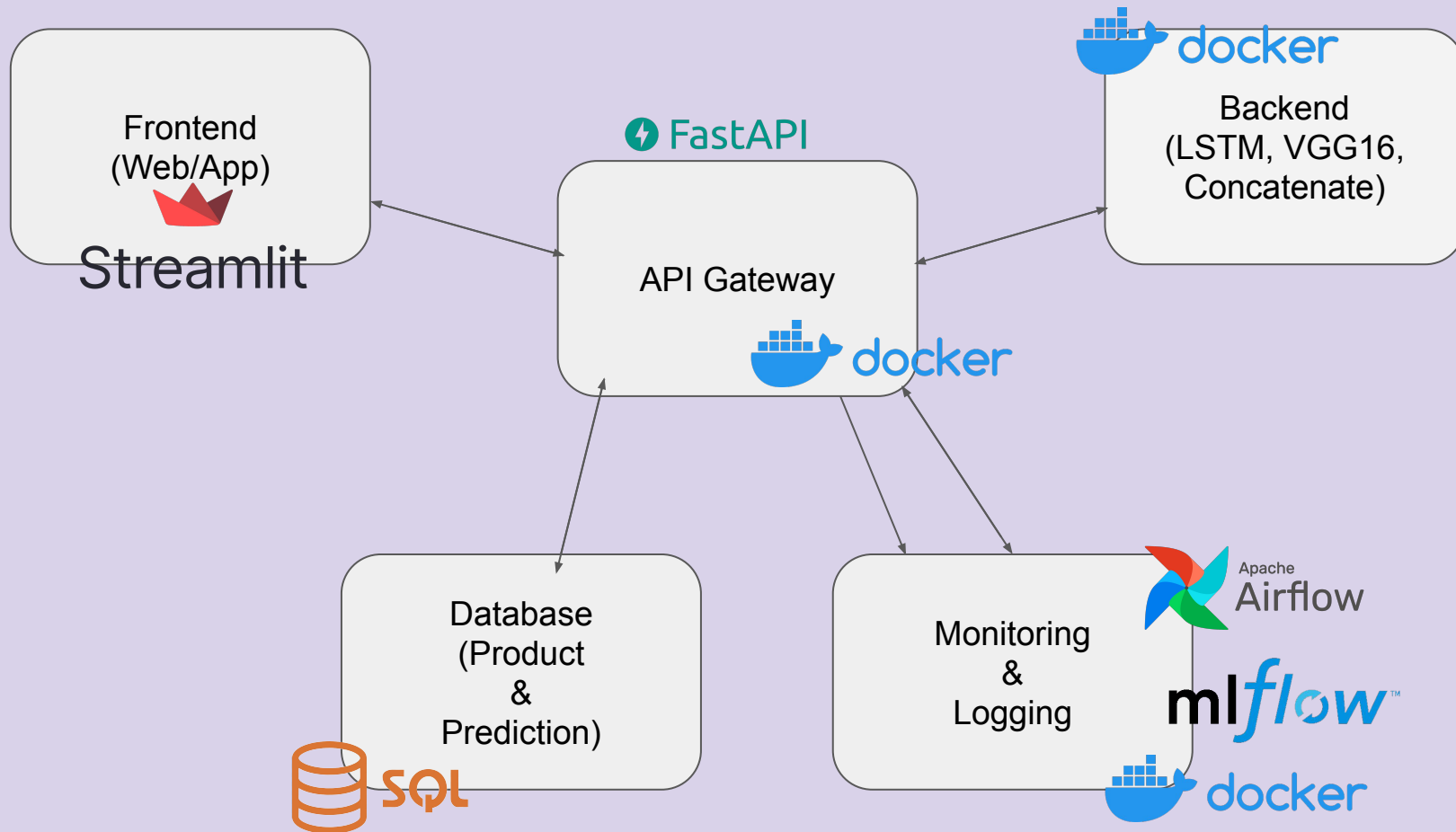
- Text LSTM model
- Image VGG16 model
- Combined model integrating outputs of the first two models
- *Purpose* of the project is to improve accuracy and efficiency for categorizing a large volume of products

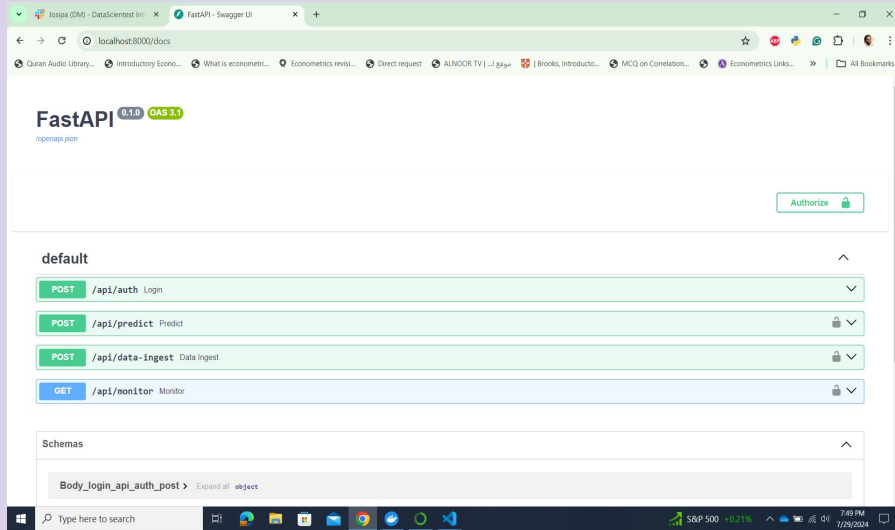
Stakeholders:

- **App Sponsor:** The e-commerce company aiming to enhance their product categorization system.
- **Users:** Internal staff, such as data analysts and product managers, who need accurate product categorization for inventory management and sales strategies.
- **Application Administrator:** IT and DevOps teams responsible for maintaining the application and ensuring its smooth operation.

Application Integration:

- **Context:** The application will be integrated into the company's existing e-commerce platform, hosted on a cloud server.
- **Medium:** The application will be accessed via an API, facilitating interaction between the models, the database, and the user interface.





- **Database Interaction:**
 - SQLite for storing product and prediction data.
 - Functions to connect, query, and manage database operations.
- **Utilities:**
 - Text and image preprocessing functions to prepare data for model input.
 - Logging for monitoring and error tracking.

Purpose: This application demonstrates integrating FastAPI with machine learning models for secure, scalable, and efficient predictions and data handling.

- **Authentication:**
 - Utilizes OAuth2 for secure authentication.
 - Access tokens generated with JWT, secured by SECRET_KEY and HS256 algorithm.
- **Machine Learning Models:**
 - **Text Model:** Best LSTM model for text classification.
 - **Image Model:** Best VGG16 model for image processing.
 - Models loaded and configured from JSON files.
- **APIs:**
 - **/api/auth:** User login, returns JWT access token.
 - **/api/predict:** Predicts using text and image inputs. Saves predictions to the database.
 - **/api/data-ingest:** Admin-only endpoint for ingesting new product data.
 - **/api/monitor:** Returns model health metrics and status.

Predictions Table:

- **Purpose:** This table stores the results of product classification predictions made by the machine learning models.
- **Role:** It acts as a repository for keeping track of prediction outputs, which can be used for analysis and further processing.

Products Table:

- **Purpose:** This table holds the detailed information about the products, such as their unique identifiers and descriptions.
- **Role:** It serves as the main source of product data that is used for generating predictions and managing inventory.

SQLite Sequence Table:

- **Purpose:** This table is part of the SQLite database system and is used to keep track of the primary keys for other tables.
- **Role:** It ensures that unique identifiers are correctly increased, supporting the integrity and consistency of the database.

Create Table

Create Index

Modify Table

Delete Table

Print

Refresh

Name	Type	Schema
<div> <div>Tables (3)</div> <div> <div> <div>></div> <div>predictions</div> <div>CREATE TABLE predictions (id INTEGER PRIMARY KEY AUTOINCREMENT, product_id INTEGER, description TEXT, in</div> </div> <div> <div>></div> <div>products</div> <div>CREATE TABLE products (id INTEGER PRIMARY KEY AUTOINCREMENT, product_id INTEGER, designation TEXT, des</div> </div> <div> <div>></div> <div>sqlite_sequence</div> <div>CREATE TABLE sqlite_sequence(name,seq)</div> </div> </div> </div>		
Indices (0)		
Views (0)		
Triggers (0)		

Table:

predictions

<

Predictions Table

- **Purpose:** Stores prediction results related to products, often used in machine learning contexts.
- **Columns:**
 - **Results Storage:** This could include various columns, such as a **Product ID** to associate predictions with specific products, predicted category or label, prediction scores, timestamps, etc.
- **Importance:** The **Predictions** table is crucial for applications that utilize machine learning to make predictions about products, such as recommending categories or predicting sales trends. Storing predictions allows for analysis, reporting, and improving predictive models.

Table:

products

Products Table

- **Purpose:** Stores information about individual products.
- **Columns:**
 - **IDs:** A unique identifier for each product. This is typically the primary key of the table, ensuring that each product can be uniquely identified.
 - **Descriptions:** Textual information about the product, which could include details such as name, features, or specifications.
 - **Image Paths:** File paths or URLs that point to images of the products. This allows the application to display product images as needed.
- **Importance:** The **Products** table is essential for storing all the key information about each product. It allows quick access to product details and images, which is vital for applications such as e-commerce platforms, where users need to view product descriptions and images.

Model Comparison Table

Model	Type	Function
Text LSTM	LSTM Neural Network	Classifies product descriptions
Image VGG16	CNN (VGG16)	Classifies product images
Concatenate	Combined Output Model	Integrates text and image data

Text LSTM Model:

- **Type:** Long Short-Term Memory (LSTM) neural network.
- **Function:** Processes and classifies product descriptions.
- **Performance:** Evaluated based on accuracy, robustness, and prediction time.

```
class TextLSTMModel:
    def __init__(self, max_words=10000, max_sequence_length=10):
        self.max_words = max_words
        self.max_sequence_length = max_sequence_length
        self.tokenizer = Tokenizer(num_words=max_words, oov_token='')
        self.model = None

    def preprocess_and_fit(self, X_train, y_train, X_val, y_val):
        self.tokenizer.fit_on_texts(X_train['description'])

        tokenizer_config = self.tokenizer.to_json()
        with open('/content/tokenizer_config.json', 'w', encoding='utf-8') as json_file:
            json_file.write(tokenizer_config)

        train_sequences = self.tokenizer.texts_to_sequences(X_train['description'])
        train_padded_sequences = pad_sequences(train_sequences, maxlen=self.max_sequence_length, padding='post',

        val_sequences = self.tokenizer.texts_to_sequences(X_val['description'])
        val_padded_sequences = pad_sequences(val_sequences, maxlen=self.max_sequence_length, padding='post', trun
```

```

class ImageVGG16Model:
    def __init__(self):
        self.model = None

    def preprocess_and_fit(self, X_train, y_train, X_val, y_val):
        # Paramètres
        batch_size = 32
        num_classes = 27

        df_train = pd.concat([X_train, y_train.astype(str)], axis=1)
        df_val = pd.concat([X_val, y_val.astype(str)], axis=1)

        # Créer un générateur d'images pour le set d'entraînement
        train_datagen = ImageDataGenerator() # Normalisation des valeurs de pixel
        train_generator = train_datagen.flow_from_dataframe(
            dataframe=df_train,
            x_col='image_path',
            y_col='prdtypecode',
            target_size=(224, 224), # Adapter à la taille d'entrée de VGG16
            batch_size=batch_size,
            class_mode='categorical', # Utilisez 'categorical' pour les entiers encodés en one-hot
            shuffle=True
        )

        # Créer un générateur d'images pour le set de validation
        val_datagen = ImageDataGenerator() # Normalisation des valeurs de pixel
        val_generator = val_datagen.flow_from_dataframe(
            dataframe=df_val,
            x_col='image_path',
            y_col='prdtypecode',
            target_size=(224, 224),
            batch_size=batch_size,
            class_mode='categorical',
            shuffle=False # Pas de mélange pour le set de validation
        )

```

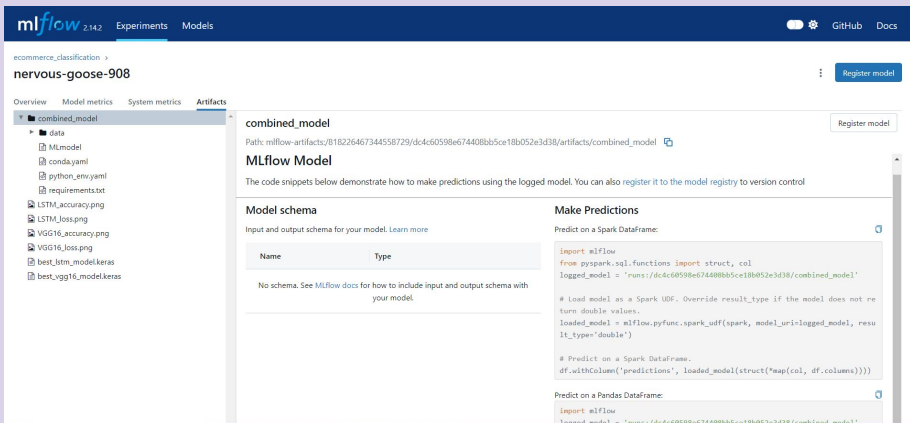
Image VGG16 Model:

- **Type:**
Convolutional Neural Network (CNN) based on the VGG16 architecture.
- **Function:**
Processes and classifies product images.
- **Performance:**
Evaluated on accuracy, robustness, and prediction time.

Concatenate Model:

- **Type:** Combines the probabilities from the Text LSTM and Image VGG16 models.
- **Function:** Integrates text and image data for final classification.
- **Performance:** Combines the strengths of the individual models for improved accuracy.

```
def optimize(self, lstm_proba, vgg16_proba, y_train):  
    # Recherche des poids optimaux en utilisant la validation croisée  
    best_weights = None  
    best_accuracy = 0.0  
  
    for lstm_weight in np.linspace(0, 1, 101): # Essayer différents poids pour LSTM  
        vgg16_weight = 1.0 - lstm_weight # Le poids total doit être égal à 1  
  
        combined_predictions = (lstm_weight * lstm_proba) + (vgg16_weight * vgg16_proba)  
        final_predictions = np.argmax(combined_predictions, axis=1)  
        accuracy = accuracy_score(y_train, final_predictions)  
  
        if accuracy > best_accuracy:  
            best_accuracy = accuracy  
            best_weights = (lstm_weight, vgg16_weight)  
  
    return best_weights
```



4. Logging and Visualization:

- **Metrics and Artifacts:** Logs accuracy, loss plots, and trained models.
- **Visualization:** Generates and saves accuracy and loss plots.

5. Model Saving:

- **Paths:** Saves models to specified file paths for future use.

6. Metrics Logging:

- **Performance Metrics:** Logs model accuracies (LSTM: 0.92, VGG16: 0.85).

Benefits:

- **Reproducibility:** Ensures experiments can be repeated with consistent results.
- **Tracking and Management:** Centralizes experiment tracking and artifact storage.
- **Automation:** Streamlines the entire ML pipeline from data ingestion to model deployment.

Key Components:

1. MLflow Setup:

racking URI: Connects to the MLflow server (<http://localhost:5000>).

Experiment: Logs data under `ecommerce_classification`.

2. Data Preprocessing:

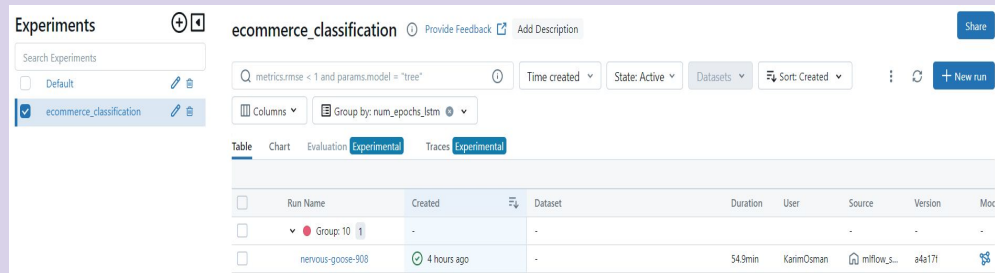
Text and Image Processing: Uses `TextPreprocessor` and `ImagePreprocessor` to clean and prepare data.

3. Model Training:

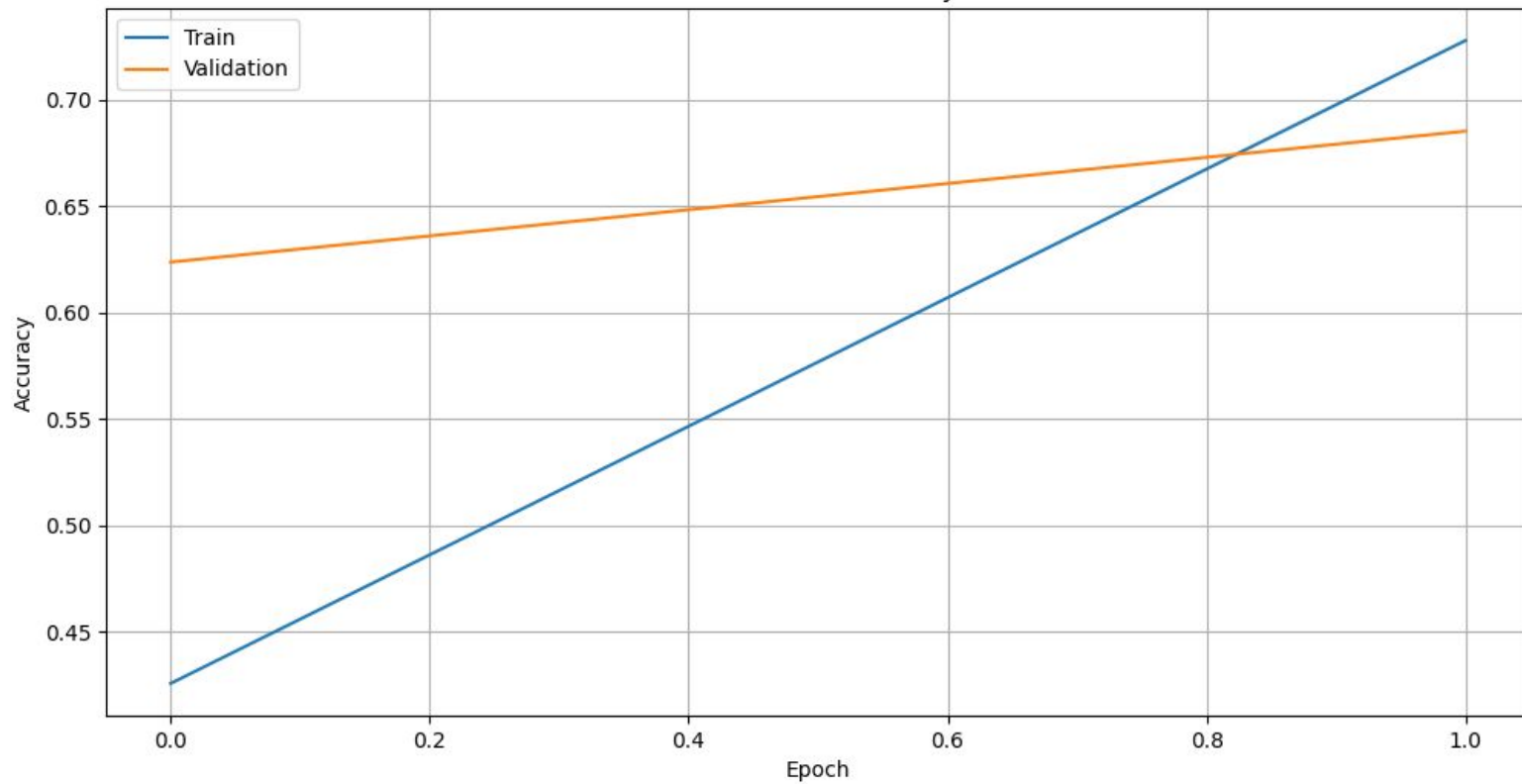
Text LSTM Model: Trains on text data.

Image VGG16 Model: Trains on image data.

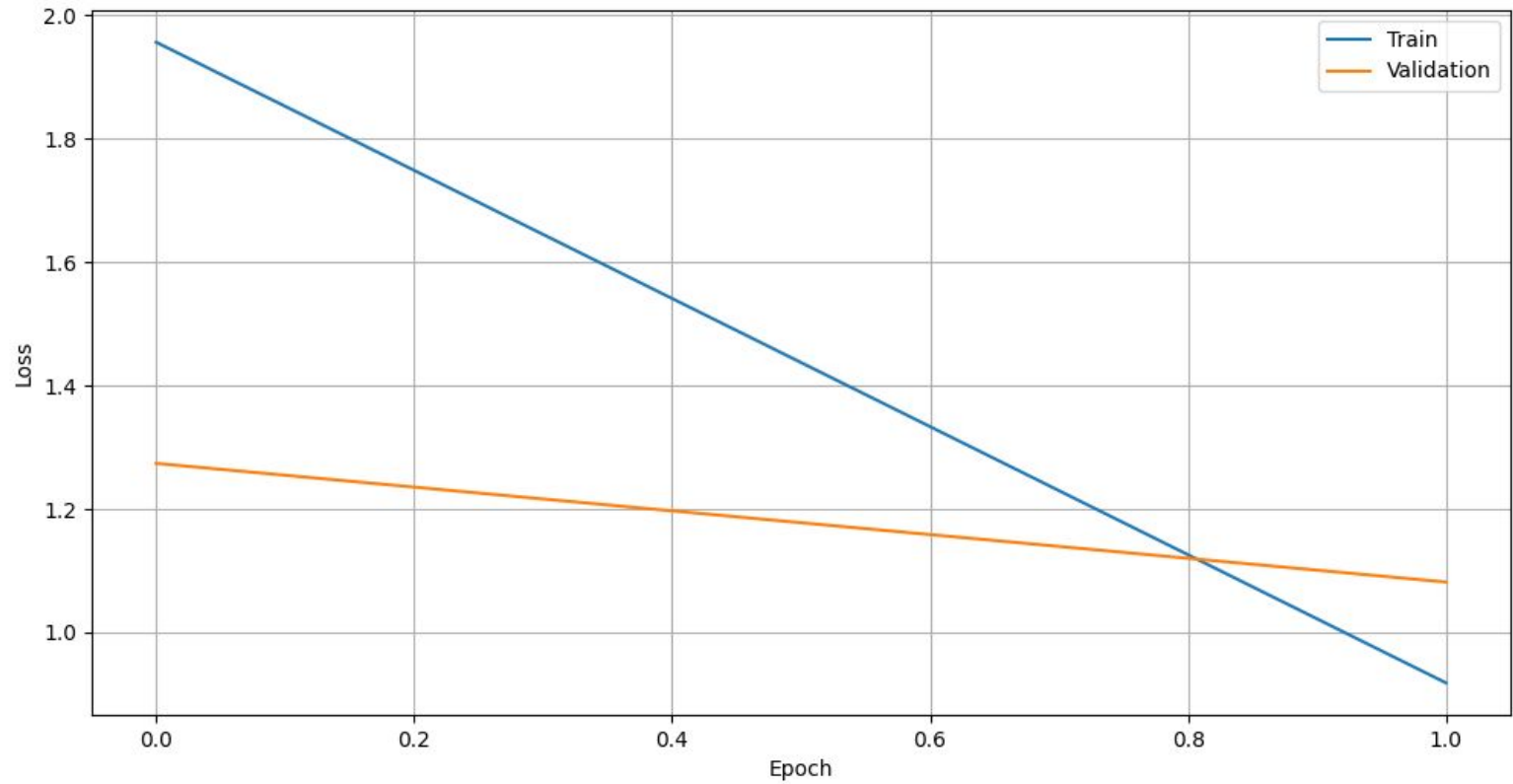
Combined Model: Integrates both LSTM and VGG16 outputs.



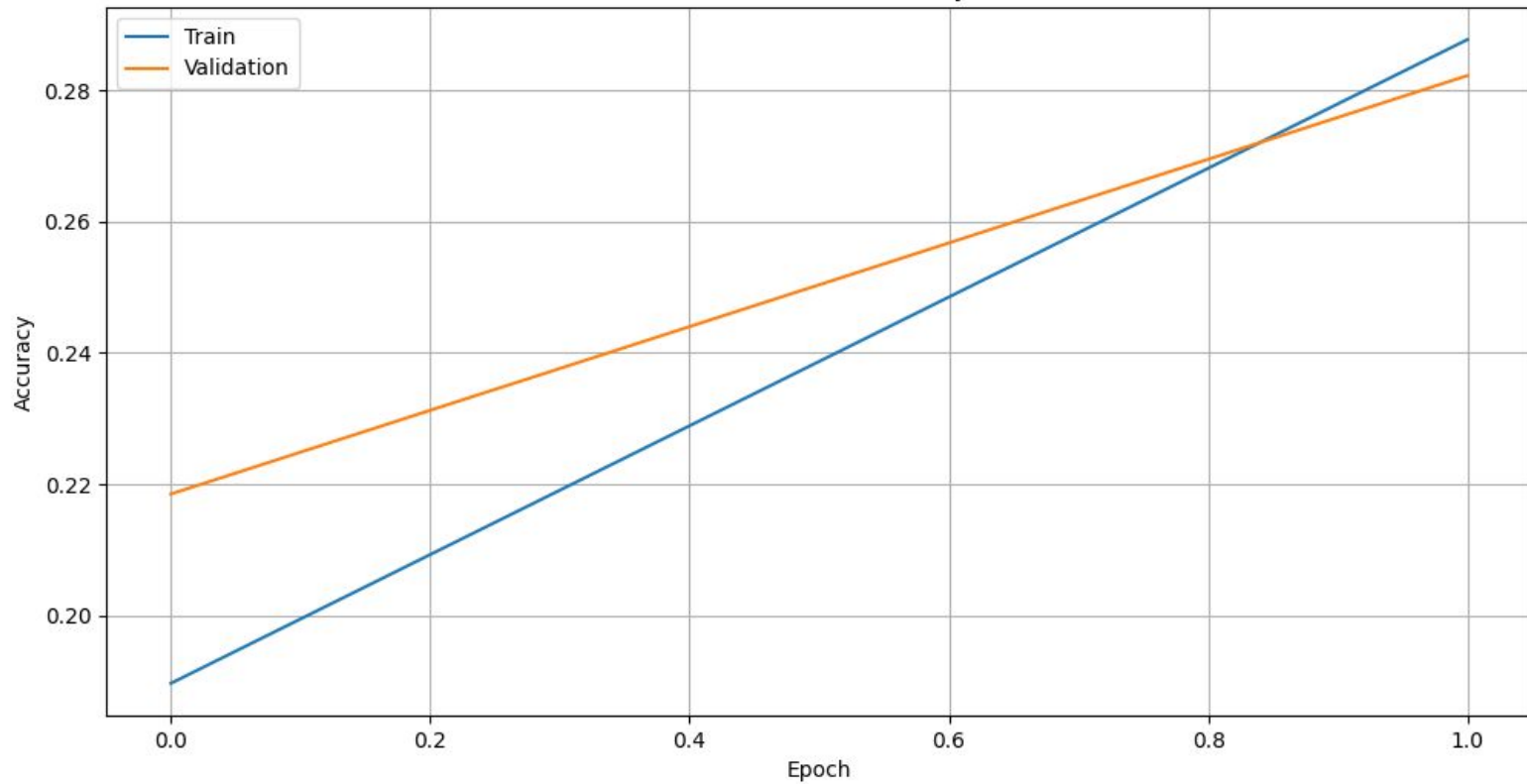
LSTM Model Accuracy



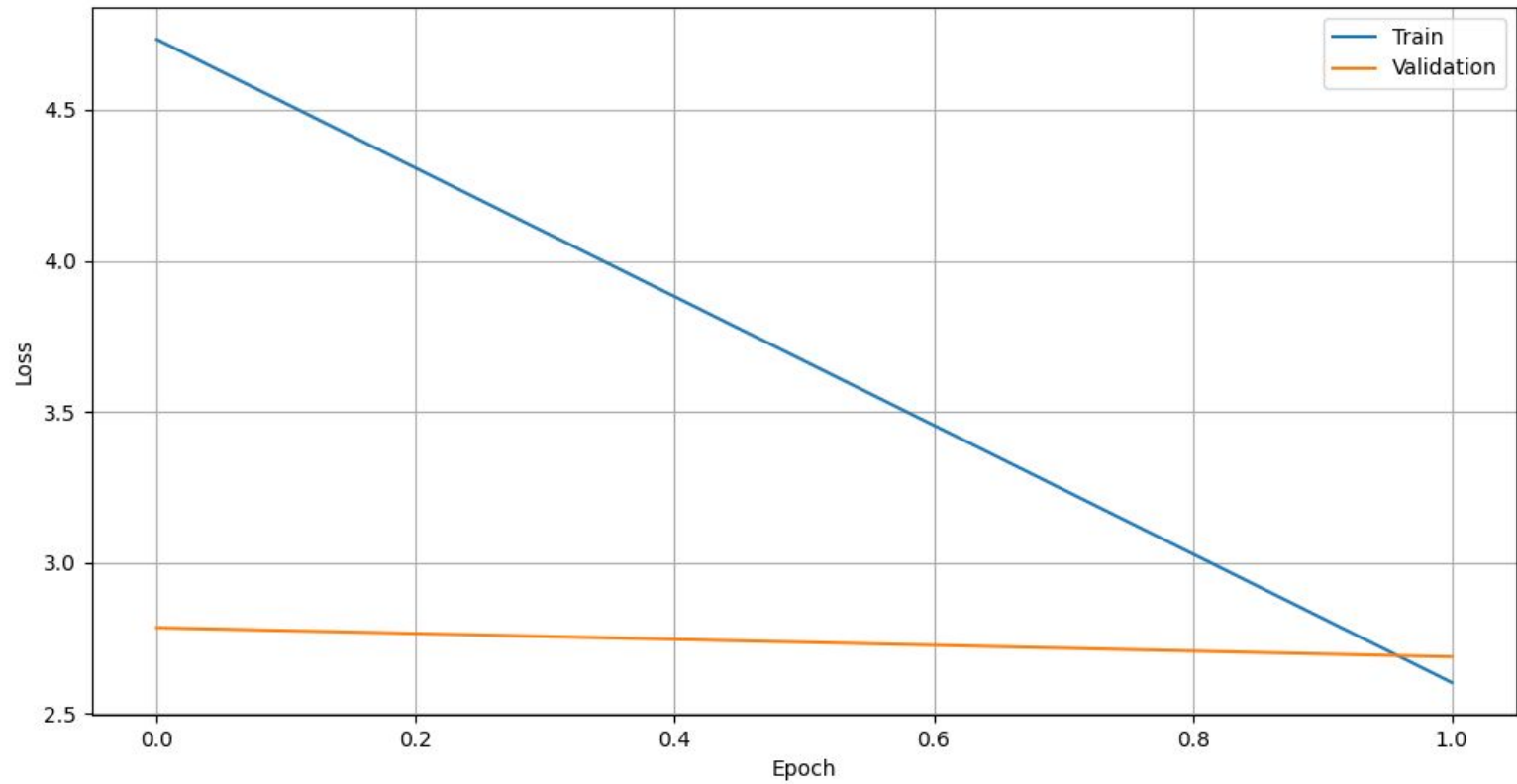
LSTM Model Loss

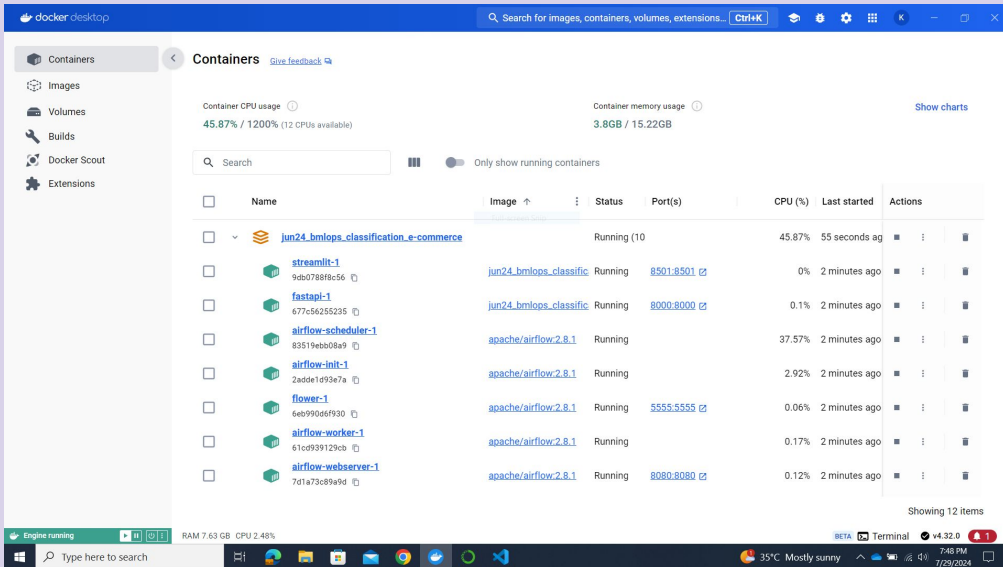


VGG16 Model Accuracy



VGG16 Model Loss





4. Portability:

- **Description:** Docker containers can run on any system that supports Docker, from local machines to cloud servers.
- **Benefit:** This portability simplifies the transition between development, testing, and production environments.

5. Simplified Deployment:

- **Description:** Dockerfiles provide a blueprint for building containers, automating the setup of environments for each component.
- **Benefit:** This automation reduces manual configuration efforts and speeds up the deployment process.

Consistent Environment:

- **Description:** Docker ensures that all components (frontend, API gateway, backend models) run in identical environments across different stages of development, testing, and production.
- **Benefit:** This consistency eliminates issues caused by differences in environments, making the deployment process more reliable.

2. Isolation and Security:

- **Description:** Each Docker container operates independently, providing isolation between different components.
- **Benefit:** This isolation enhances security and ensures that changes in one component do not affect others.

3. Scalability:


- **Description:** Docker makes it easy to scale individual components by running multiple instances of a container as needed.
- **Benefit:** This flexibility supports the demands of varying workloads, especially important for handling large volumes of data in real-time.

Pipelines:

- **Automation of Testing, Building, and Deploying:**
 - **Description:** CI/CD pipelines automate the process of testing code, building Docker images, and deploying them to production environments. This ensures that code changes are integrated and deployed smoothly.
 - **Benefit:** Reduces manual intervention, accelerates deployment cycles, and improves the reliability of deployments.
- **Workflows Based on ci.yml and deployment.yml:**
 - **ci.yml:** Defines the steps for continuous integration, including code checkout, dependency installation, and testing.
 - **deployment.yml:** Specifies the deployment process, including building Docker images and deploying them to the target environment.

```
1 name: AIRFLOW WORKFLOW
2
3 on:
4   push:
5     branches:
6       - main
7
8 jobs:
9   airflow-setup:
10    runs-on: ubuntu-latest
11
12    steps:
13      - name: Checkout code
14        uses: actions/checkout@v3
15
16      - name: Set up Python
17        uses: actions/setup-python@v4
18        with:
19          python-version: 3.9
20
21      - name: Install dependencies
22        run: |
23          pip install -r airflow/requirements.txt
24
25      - name: Build Docker image
26        run: docker build -t airflow-image -f airflow/Dockerfile .
27
28      - name: Run Airflow Docker container
29        run: |
30          docker run -d -p 8080:8080 airflow-image
31
32      - name: Check Airflow webserver status
33        run: curl -sSf http://localhost:8080
```



Deploy 

Login

Username

Password



Login

Model Prediction

Upload an Image



Drag and drop file here

Limit 200MB per file • JPG, JPEG, PNG

Browse files

Product Description

