# CS5830 Big Data Lab Project

Group-2
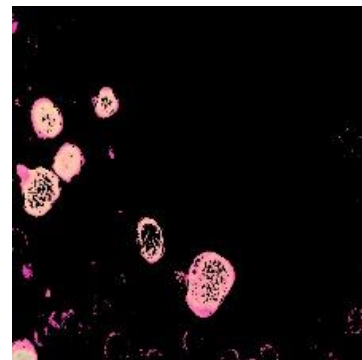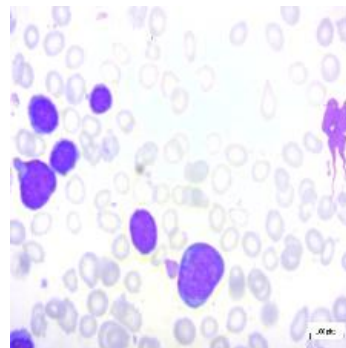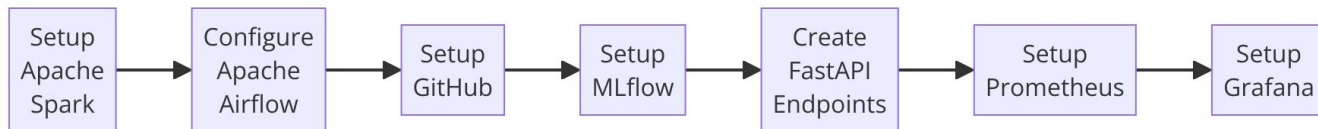
Vishal V ME20B204
Akranth Reddy ME20B100
Sai Gowtham ED19B063

# Problem Statement & Workflow

Acute Lymphoblastic Leukemia (ALL) is an aggressive form of cancer that predominantly affects children. Early detection and accurate classification of ALL are crucial for effective treatment and improving patient outcomes. However, diagnosing ALL through microscopic blood smear images is challenging due to visual similarities with other conditions and the need for specialized expertise.

This project aims to leverage an MLOps approach to build an end-to-end machine learning solution for the detection and classification of ALL from microscopic blood smear images. The solution will include a

- data preprocessing pipeline using Apache Airflow,
- machine learning model tracking via MLflow, and
- a scalable REST API for model deployment using FastAPI.
- the entire solution will be containerized for seamless deployment and monitored using Prometheus and Grafana.



```
Setup Apache Spark → Configure Apache Airflow → Setup GitHub → Setup MLflow → Create FastAPI Endpoints → Setup Prometheus → Setup Grafana
```

# Data preprocessing and training

# Apache Spark

## Without Spark

```python
def preprocessing(car_path, mask_path):
    car_img = tf.io.read_file(car_path)
    car_img = tf.image.decode_jpeg(car_img, channels=3)
    car_img = tf.image.resize(car_img, img_size)
    car_img = tf.cast(car_img, tf.float32) / 255.0

    mask_img = tf.io.read_file(mask_path)
    mask_img = tf.image.decode_jpeg(mask_img, channels=3)
    mask_img = tf.image.resize(mask_img, img_size)
    mask_img = mask_img[:,:,:1]
    mask_img = tf.math.sign(mask_img)

    return car_img, mask_img
```

## With Spark

```python
def preprocessing(car_path, mask_path):
    # Read car image and mask image
    car_img_data = tf.io.read_file(car_path)
    mask_img_data = tf.io.read_file(mask_path)

    # Define a function to process each image
    def process_image(img_data):
        img = tf.image.decode_jpeg(img_data, channels=3)
        img = tf.image.resize(img, img_size)
        img = tf.cast(img, tf.float32) / 255.0
        return img.numpy()

    # Process car image and mask image
    car_img_np = np.array([process_image(car_img_data)])
    mask_img_np = np.array([process_image(mask_img_data)[:,:,:1]

    return car_img_np, mask_img_np

# Define the image size
img_size = (256, 256)

# Call the preprocessing function
car_img_np, mask_img_np = preprocessing(car_path, mask_path)

# Create RDDs from the numpy arrays
car_img_rdd = spark.sparkContext.parallelize(car_img_np)
mask_img_rdd = spark.sparkContext.parallelize(mask_img_np)

# Collect RDDs into lists
car_img_list = car_img_rdd.collect()
mask_img_list = mask_img_rdd.collect()

# Close the SparkSession
spark.stop()

# Convert lists to numpy arrays
car_img_np_final = np.array(car_img_list)
mask_img_np_final = np.array(mask_img_list)
```
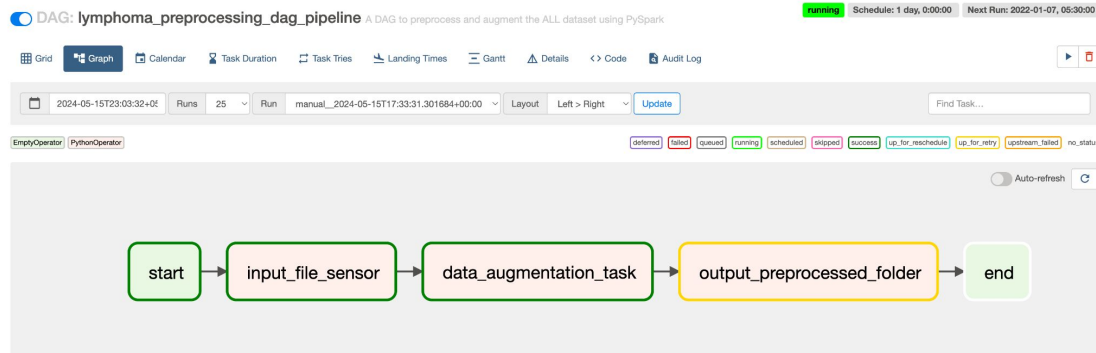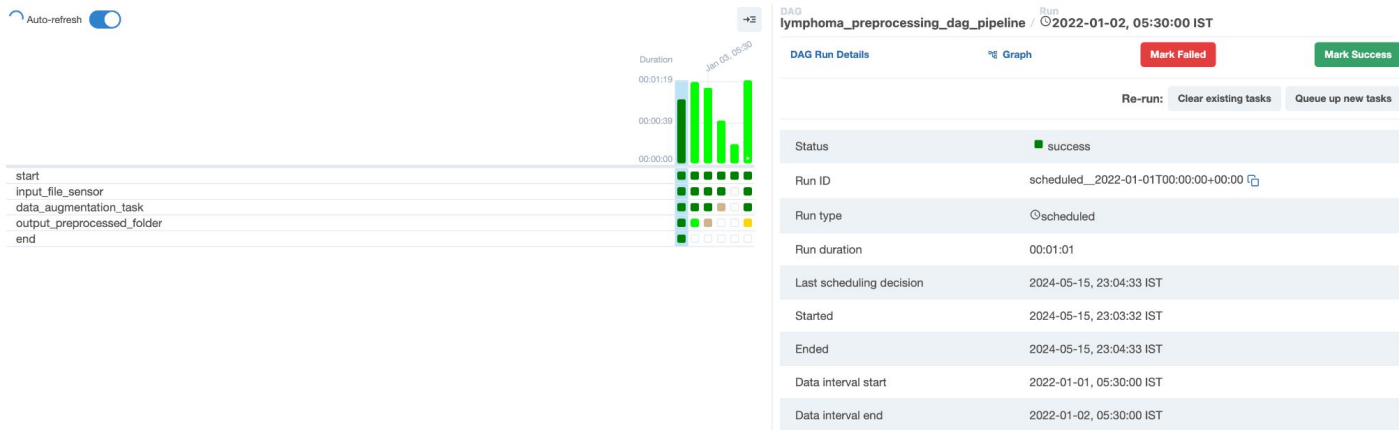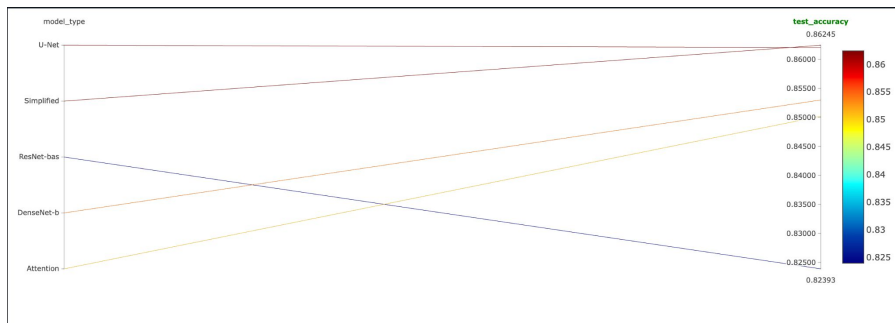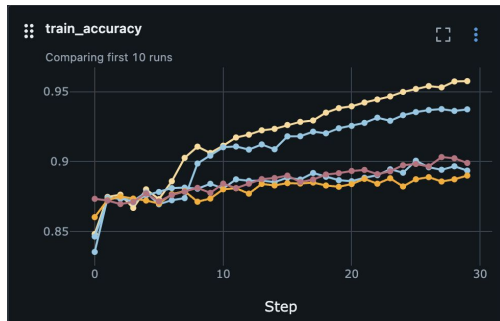
# Airflow



## Use

- Automation: Reduces manual effort and ensures consistent preprocessing and augmentation of images.
- Reproducibility: Provides a repeatable process that can be easily triggered and monitored.
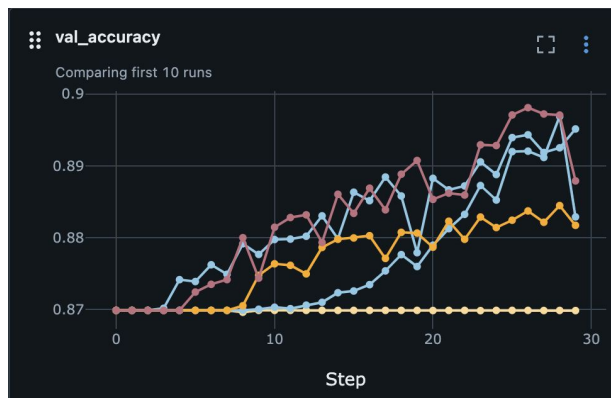- Flexibility: Can be adapted for different datasets and preprocessing requirements.

# MLflow



- The experiment tracking with MLflow provides valuable insights into the performance of different models
- DenseNet-based Model: Offers a balanced performance with a high training accuracy (93.73%) and a strong validation accuracy (88.29%). Its consistent metrics make it a reliable choice for further use.

# Deployment - the project

| | | |
|---|---|---|
| 📁 data | initial commit | 20 hours ago |
| 📁 src | prometheous is working | 32 minutes ago |
| 📁 utils | initial commit | 20 hours ago |
| 📁 weights | initial commit | 20 hours ago |
| 📄 .dockerignore | update 16th may evening | 18 hours ago |
| 📄 .gitattributes | initial commit | 20 hours ago |
| 📄 .gitignore | initial commit | 20 hours ago |
| 📄 Dockerfile | commit early morning | 8 hours ago |
| 📄 README.md | initial commit | 20 hours ago |
| 📄 docker-compose.yml | prometheous is working | 32 minutes ago |
| 📄 fastapi.json | prometheous is working | 32 minutes ago |
| 📄 prometheus.yml | prometheous is working | 32 minutes ago |
| 📄 requirements.txt | prometheous is working | 32 minutes ago |
| 📄 task1.py | prometheous is working | 32 minutes ago |

# Fast-api

```python
35    # Define counters for API usage
36    api_usage_counter = Counter("api_usage", "API usage counter", ["client_ip"])
37
38    Instrumentator().instrument(app).expose(app)
39
40    # Define gauges for API processing time
41    api_time_gauge = Gauge("api_time", "API processing time", ["client_ip"])
42    api_time_per_char_gauge = Gauge("api_time_per_char", "API processing time per character", ["client_ip"])
43
44  > def load_model(): …
48
49  > def preprocess_image(image): …
55
56  > def predict_label(image, model): …
63
64
65    @app.post("/predict/")
66  > async def predict(request:Request, file: UploadFile = File(...)): …
106
107   if __name__ == "__main__":
108      uvicorn.run(app, host="127.0.0.0", port=8080)
109
```

# Swagger-ui, client side

Input image

file ★ required
string($binary)    [ Browse... ] WBC-Malignant-Pro-001.jpg

| Execute | Clear |

## Responses

Curl

```
curl -X 'POST' \
  'http://localhost:8080/predict/' \
  -H 'accept: application/json' \
  -H 'Content-Type: multipart/form-data' \
  -F 'file=@WBC-Malignant-Pro-001.jpg;type=image/jpeg'
```

Request URL

```
http://localhost:8080/predict/
```

Server response

Code        Details

200
            Response body



Output predicted masked image.

# Prometheus

# Graphana



**Metric 1:** Number of Requests: Tracks the total number of requests received by the FastAPI server.ions.

**Metric 2:** CPU Usage Percentage: Monitors the percentage of CPU utilization by the FastAPI application.

**Metric 3:** Network I/O (Sent and Received Bytes): Measures the total bytes sent by the server, reflecting the volume of data served to the clients.

**Metric 4:** Memory Usage: Tracks the memory usage in bytes of the FastAPI application.

# Docker

# Git

- Our project had lots of training data, ~190 MB.
- Used Git - LFS for tracking changes in data
- And the whole project  is tracked using Git version control.
- Github link:

# Contributions

Version control: **Git** and **Git-lfs**

Vishal V (ME20B204):

> **Apache Airflow** Preprocessing Pipeline, Experiments Tracking using **MLFlow**, Monitoring Dashboard using **Grafana**

Akranth (ME20B100):

> Deployed APIs with **FastAPI**, tracking the API usage with **prometheus** and visualization with **grafana**, containerized the whole thing using **docker**. Tracking the project using git.

Sai Gowtham Tamminaina (ED19B063)

> Apache Spark and Airflow for the Preprocessing Pipeline