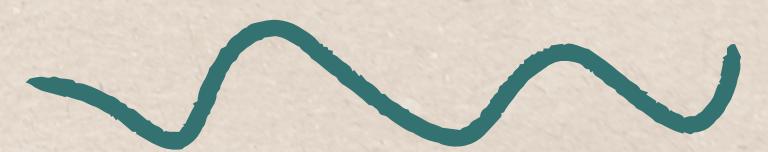




AI-Driven Prediction of Kubernetes Cluster Failures

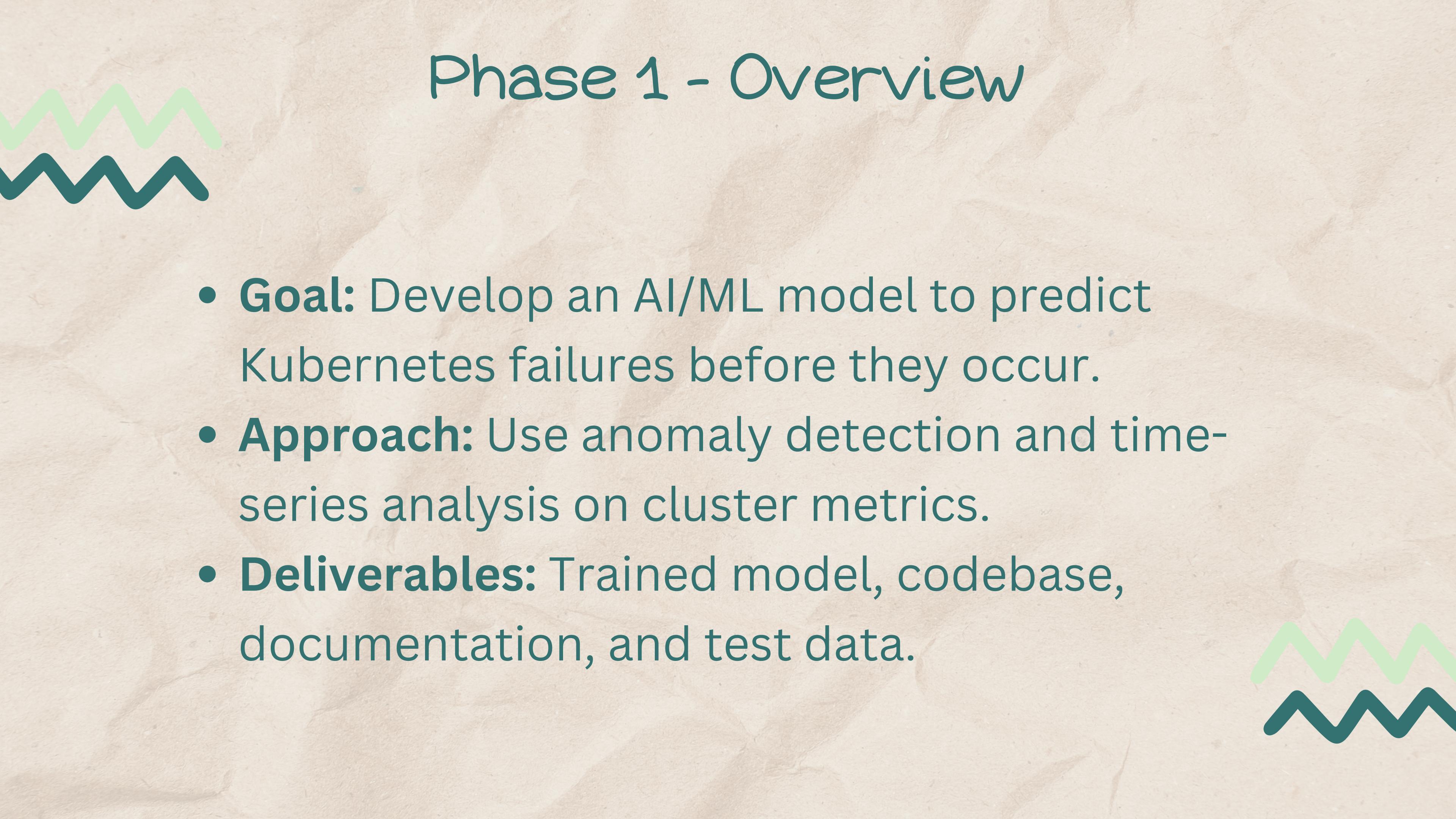


**Phase 1 - Prediction Model
Presented By Capital W**



Introduction

Kubernetes clusters experience failures such as pod crashes, resource exhaustion, and network issues. These failures can lead to downtime, increased costs, and reduced performance. The challenge is to build an AI/ML model that can accurately predict these failures before they occur by analyzing historical and real-time cluster metrics.



Phase 1 - Overview

- **Goal:** Develop an AI/ML model to predict Kubernetes failures before they occur.
- **Approach:** Use anomaly detection and time-series analysis on cluster metrics.
- **Deliverables:** Trained model, codebase, documentation, and test data.



Data Collection

Data Sources:

<https://drive.google.com/drive/folders/1FkeDE2eBsoe8DOU2h0dqoVpsgzOEfW9M?usp=sharing>

Metrics Considered:

- CPU & Memory Utilization
- Pod/Node Status
- Network Latency & Throughput
- Disk I/O & Logs





Data Collection

Resource Utilization

- CPU Usage (%)
 - Memory Usage (%), MB
 - Network I/O (Bytes, Packets)
 - Filesystem Operations (Reads/Writes)
- 

Data Collection

Pod Status

- Current Status (e.g., Running)
- Restarts
- Ready/Total Containers



Data Collection

Events

- Pod Events (Type, Reason, Age, Message)
 - Node Events (Reason, Age, Message)
- 



Data Collection

Resource Allocation

- Memory Requests/Limits
 - Memory Utilization vs Limits
- 



MODEL ARCHITECTURE

ML Techniques Used

- LSTM-Based Time-Series Forecasting for anomaly detection.
- Anomaly Classification to categorize failure types (CPU overload, memory exhaustion, etc.).
- Confidence Scoring Mechanism for accurate predictions and early warnings.

LIBRARIES & TOOLS

- **Deep Learning:** TensorFlow, PyTorch for LSTM model training.
- **Traditional ML:** Scikit-learn for feature engineering & preprocessing.
- Kubernetes real-time metrics collection.
- Prometheus & Grafana for visualization and monitoring.

KEY IMPLEMENTATIONS

- Rolling Window Statistics & Rate-of-Change Features for better anomaly detection.
- Adaptive Thresholding Mechanism to refine anomaly alerts dynamically.
- Real-Time Model Deployment Strategy for integration with Kubernetes clusters.

Results #1

LSTM-Based Model For Prediction

Trained an LSTM-based model to predict anomalies using historical and real-time metrics.

DEMO:

https://drive.google.com/file/d/1_d42KJCyPbpIZhFdiaoI99gHWLneAl2/view?usp=drive_link

The screenshot shows a code editor interface with several tabs open. The main tab displays a Python script named `app1.py` containing the following code:

```
File Edit Selection View Go Run Terminal Help
EXPLORER ... inergering.py 1.M modevaluation.py 1.A earlywarning.py 1.A app1.py k8integration.py A Containerization A basicimp.py A container2 A
K8
> __pycache__
> .venv
> requirements.txt
adapmoni.py A
anamolyprediction.py 1.M
app.py A
app1.py 1.A
basicimp.py A
container2 A
Containerization A
datapreprocessing.py A
dataSynthetic.csv A
Dockerfile A
earlywarning.py 1.A
errortate.py A
featureengineering.py A
fetch_metrics.py A
k8integration.py A
k8s_failure_model.pkl !
k8s-deployment.yaml !
lstm_anomaly_model.h5 M
lstmmodel.py A
lstmsequence.py A
modevaluation.py A
multivariate.py A
patch.json A
prometheus_data.csv A
prometheus_data.json A
threshold.py A
train_model.py A
weightedloss.py A
PROBLEMS 5 OUTPUT DEBUG CONSOLE TERMINAL PORTS
[2025-03-25 22:43:50.601139] No anomalies detected in current window.
1/1 0s 21ms/step
[2025-03-25 22:43:51.652088] No anomalies detected in current window.
1/1 0s 22ms/step
[2025-03-25 22:43:52.701210] No anomalies detected in current window.
1/1 0s 21ms/step
[2025-03-25 22:43:53.751313] WARNINGS DETECTED:
- opentelemetry-demo-shippingservice-86ccddbd5b-9hjhw: Unknown (Probability: 0.74)
- opentelemetry-demo-cartservice-56c6b98587-tsfc: Memory Exhaustion Risk (Probability: 0.92)
- opentelemetry-demo-paymentservice-5fcfd5559c-kx9hl: Unknown (Probability: 0.92)
1/1 0s 22ms/step
[2025-03-25 22:43:54.803380] WARNINGS DETECTED:
- opentelemetry-demo-paymentservice-5fcfd5559c-kx9hl: Unknown (Probability: 0.92)
- opentelemetry-demo-ffspostgres-7464b97d47-qz9zg: Unknown (Probability: 0.85)
1/1 0s 21ms/step
[2025-03-25 22:43:55.854046] No anomalies detected in current window.
1/1 0s 21ms/step
> OUTLINE
> TIMELINE
main+ 0 0 △ 5
```

The terminal pane at the bottom shows the execution of the script, with the last few lines indicating no anomalies detected and some warning messages detected.

Results #2

Real-time Kubernetes Data Collection

```
PS C:\Kube\W\k8s-failure-dataset-generator> python3 dataset-generator1.py
Fetching data for pods in namespace monitoring
```

- we attempted to collect metrics from a Kubernetes cluster (ran on a local machine) using Minikube .
- Minikube is a tool that allows you to run a single-node Kubernetes cluster locally for development and testing purposes.

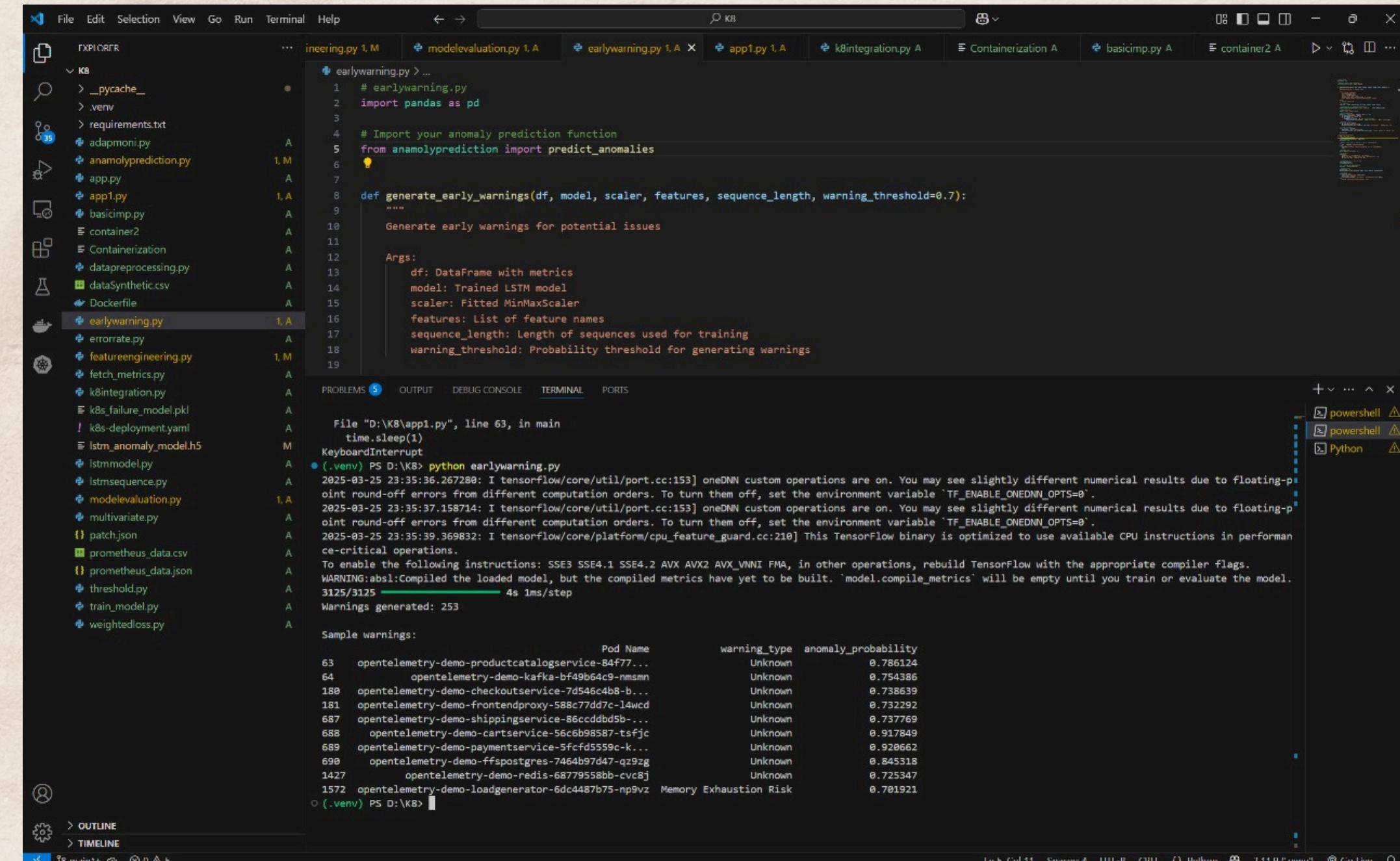
Results #3

Alert System For Anomaly Detection

Developed a proactive alert system that categorizes warnings based on thresholds and generates actionable insights.

DEMO:

https://drive.google.com/file/d/1_d42KJCyPbplZhF-diaoI99gHWLneAl2/view?usp=drive_link



The screenshot shows a code editor interface with several files open in tabs. The main file is `earlywarning.py`, which contains the following code:

```
# earlywarning.py > ...
# Import your anomaly prediction function
from anomalyprediction import predict_anomalies

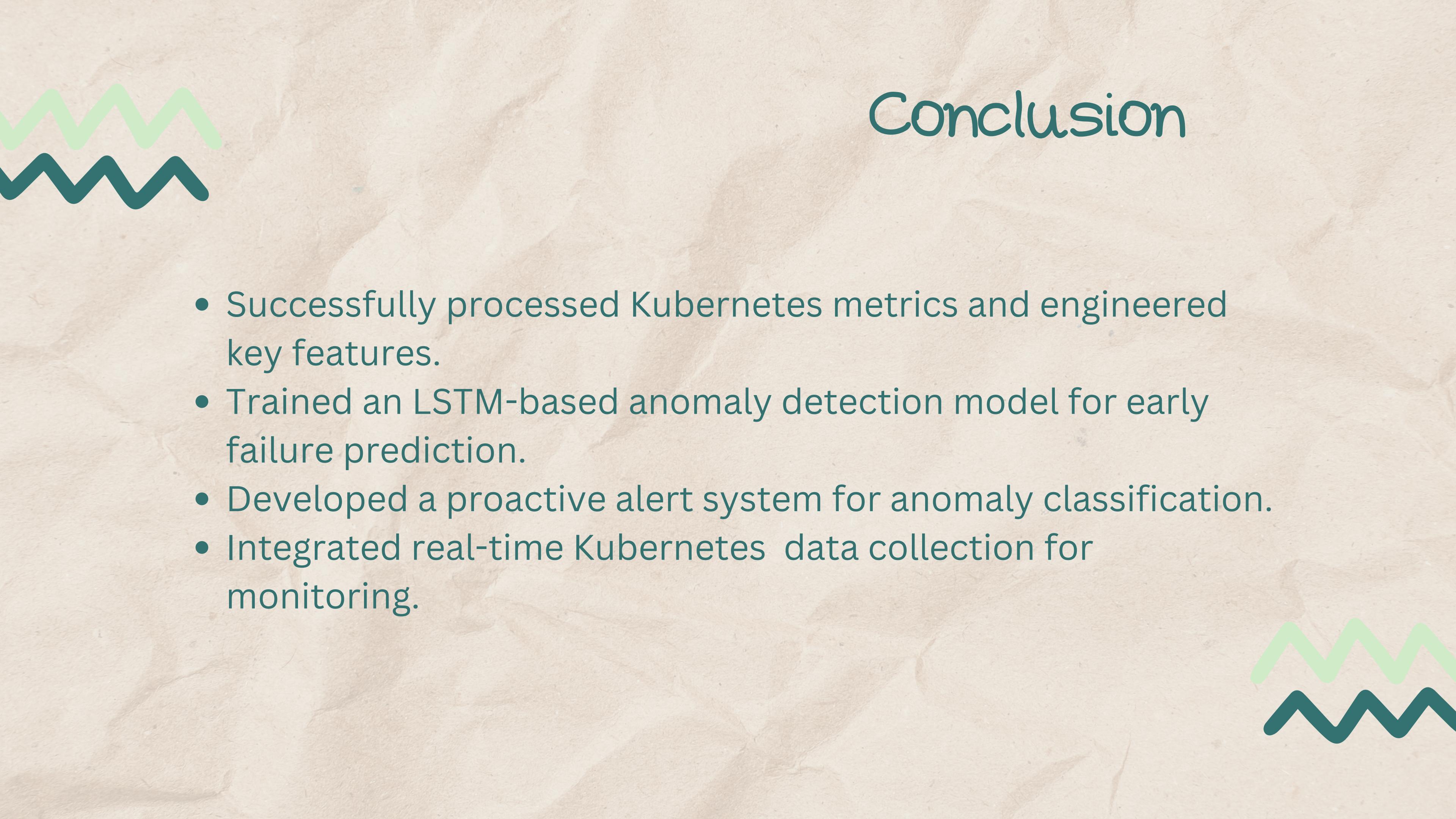
def generate_early_warnings(df, model, scaler, features, sequence_length, warning_threshold=0.7):
    """
    Generate early warnings for potential issues

    Args:
        df: DataFrame with metrics
        model: Trained LSTM model
        scaler: Fitted MinMaxScaler
        features: List of feature names
        sequence_length: Length of sequences used for training
        warning_threshold: Probability threshold for generating warnings
    """

    # Implementation of generate_early_warnings function
```

The code editor also displays a terminal window showing the execution of `earlywarning.py` and its output, which includes TensorFlow compilation logs and sample warning data. The sample warning data is as follows:

Pod Name	warning_type	anomaly_probability
63	opentelemetry-demo-productcatalogservice-84f77...	Unknown 0.786124
64	opentelemetry-demo-kafka-bf49b64c9-nmsmn	Unknown 0.754386
180	opentelemetry-demo-checkoutservice-7d546c4b8-b...	Unknown 0.738639
181	opentelemetry-demo-frontendproxy-588c77dd7e-14wcd	Unknown 0.732292
687	opentelemetry-demo-shippingservice-86ccddbd5b-...	Unknown 0.737769
688	opentelemetry-demo-cartservice-56c6b98587-tsfjc	Unknown 0.917849
689	opentelemetry-demo-paymentservice-5fcfd5559c-k...	Unknown 0.920662
690	opentelemetry-demo-ffspostgres-7464b97d47-qz9zg	Unknown 0.845318
1427	opentelemetry-demo-redis-68779558bb-cvc8j	Unknown 0.725347
1572	opentelemetry-demo-loadgenerator-6dc4487b75-np9vz	Memory Exhaustion Risk 0.701921



Conclusion

- Successfully processed Kubernetes metrics and engineered key features.
- Trained an LSTM-based anomaly detection model for early failure prediction.
- Developed a proactive alert system for anomaly classification.
- Integrated real-time Kubernetes data collection for monitoring.

Our Approach

1. LSTM-based anomaly detection for Kubernetes failures.
2. Confidence-based classification of anomalies (e.g., CPU overload, memory exhaustion).
3. Real-time integration with Kubernetes APIs for continuous monitoring.
4. Proactive alert system with adaptive thresholding.

Are We Going in the Right Direction?

1. Does our approach align with best practices for AI-driven Kubernetes monitoring?
2. Are there any gaps or areas for refinement in our methodology?
3. How can we further optimize our prediction accuracy and response time?
4. Should our Kubernetes Anomaly Detection System be developed as a service-based or product-based solution?



Thank
you