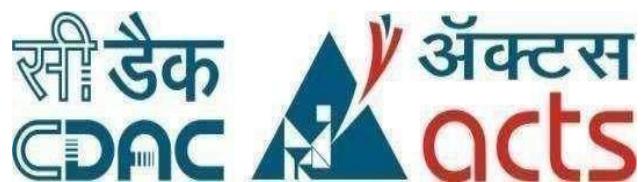


PROJECT REPORT ON

HPC Node/Cluster Failure Prediction With prevention



Submitted By

Suraj Kumar Choudhary (240840127041)

*Submitted in partial fulfillment for the award of
Graduate Diploma in High-Performance Computing
System Administration from C-DAC ACTS (Pune)*

Under the Guidance of

Roshan Gami

Abstract

High-Performance Computing (HPC) clusters are essential for executing large-scale computations across multiple nodes. However, failures in these clusters can lead to significant downtime, reduced performance, and increased maintenance costs. This project aims to develop a **failure prediction system** for HPC nodes by leveraging **log analysis, real-time monitoring, and machine learning techniques**.

The proposed system collects **system logs, job scheduler logs (e.g., SLURM, PBS), and hardware monitoring metrics (CPU, memory, disk I/O, network traffic, etc.)** from compute nodes. Using **log parsing tools (Logstash, Fluentd)** and **monitoring frameworks (Prometheus, Grafana)**, the system preprocesses the data to extract meaningful features. Machine learning models, including **anomaly detection (Isolation Forest, DBSCAN)** and **predictive modeling (XGBoost, LSTMs, RNNs)**, analyze historical logs and real-time data to predict potential node failures.

The predictive model is integrated into a **real-time monitoring and alerting system**, where failure probabilities are visualized in **Grafana dashboards** and alerts are triggered via **AlertManager (email, Slack, PagerDuty)** when critical thresholds are met. The project also addresses key challenges such as **handling large log volumes, dealing with sparse failure data, and ensuring model interpretability**.

By implementing this system, **HPC administrators can proactively detect and mitigate failures, reducing downtime, optimizing resource usage, and improving overall system reliability**. The solution is designed for **scalability** using **Docker and Kubernetes**, ensuring seamless deployment and integration with existing HPC infrastructures.

This project demonstrates the **practical application of AI/ML in HPC environments**, enhancing system resilience through **proactive failure detection and predictive maintenance**.

<u>Table of Contents</u>	Page no
1. Introduction	7-9
• Overview	
• Background	
• Importance of HPC Failure Prediction	
• Challenges in HPC Failure Prediction	
• Methods for HPC Failure Prediction	
2. Methods for HPC Failure Prediction	
• Data Collection and Preprocessing	
• Machine Learning Techniques for Failure Prediction	
• Real-Time Monitoring and Prediction System	
3. Motivation	10-11
• Why Failure Prediction is Essential?	
• Key Motivations for This Project:	
○ Minimize Downtime	
○ Enhance System Reliability	
○ Improve Resource Utilization	
○ Enable Data-Driven Decision Making	

○ Reduce Maintenance Costs	
○ Broader Implications	
4. Problem Statement	12-13
• Overview of Challenges in HPC Environments	
• Key Challenges Addressed in This Project	
5. Objectives	14-17
• Proactive Maintenance	
• Log Collection and Analysis	
• Real-Time Monitoring	
• Machine Learning-Based Failure Prediction	
• Integration with Monitoring Systems	
• Scalable Deployment	
• Comparison of Predictive Models	
• Implementation of a Decision Support System	
• Performance Testing and Optimization	
• Deployment in a Real HPC Environment	
6. Planning	18-21
• Introduction	
• Planning Objectives	
• Major Issues Addressed in Project Plan	

• Planning Activities	
• Waterfall Model	
7. Methodology	22-25
• Methodology Overview	
• Step 1: Data Collection	
• Step 2: Data Preprocessing & Feature Engineering	
• Step 3: Exploratory Data Analysis (EDA)	
• Step 4: Machine Learning-Based Prediction	
• Workflow Summary	
8. Implementation Steps	26-31
8.1 System Requirements	
• Hardware Requirements	
• Software Requirements	
8.2 Installation and Setup	
• Setting Up the HPC Cluster	
• Installing Required Software	
• Configuring Log Collection and Monitoring	
8.3 Data Collection	
• Collecting System Logs	
• Monitoring Hardware Metrics	

• Annotating Failure Events	
8.4 Data Preprocessing	
• Log Parsing	
• Feature Engineering	
• Handling Missing Data	
8.5 Exploratory Data Analysis (EDA)	
8.6 Model Development	
8.7 Deployment	
8.8 Validation and Testing	
8.9 Using the Dashboard	
8.10 Troubleshooting	
9. Screenshots	31-61
• Step-by-Step Implementation Screenshots	
10. Future Scope	62-64
11. Conclusion	65
12 References	66-67

INTRODUCTION

OVERVIEW

1. Background

High-Performance Computing (HPC) clusters play a crucial role in scientific research, data analytics, artificial intelligence, and various computational fields that require extensive processing power. These clusters consist of multiple interconnected nodes working together to execute large-scale computations efficiently. However, node failures and performance degradation are common challenges that can significantly impact computational tasks, leading to increased downtime and maintenance efforts.

Failures in HPC environments can result from various factors, including hardware malfunctions, software crashes, overheating, network congestion, and disk failures. Predicting these failures in advance can help administrators take proactive measures to minimize disruptions and optimize resource utilization. Traditional failure detection relies on reactive maintenance, where issues are addressed only after they occur, leading to inefficiencies, prolonged downtime, and increased operational costs. In contrast, a proactive failure prediction system enables early detection and mitigation of potential failures, improving the overall reliability and efficiency of HPC clusters.

2. Importance of HPC Failure Prediction

HPC clusters are widely used in domains such as climate modeling, bioinformatics, financial simulations, and engineering design. These fields demand continuous system availability and minimal disruptions, making failure prediction an essential component of HPC management. When failures are detected too late, computational jobs may be lost, requiring reruns and leading to increased energy consumption and resource wastage. A robust failure prediction model helps in preventing such inefficiencies by identifying potential issues before they escalate.

Modern HPC systems generate vast amounts of logs and monitoring data, including system metrics, application logs, and hardware performance indicators. These logs provide valuable insights into system health and operational status. By leveraging log data, administrators can detect early warning signs of failures and take corrective actions to avoid catastrophic breakdowns. The combination of log analysis and machine learning techniques enhances the ability to predict failures accurately, reducing overall maintenance costs and improving system uptime.

3. Challenges in HPC Failure Prediction

Despite the advantages of failure prediction, several challenges exist in developing an accurate and efficient predictive model. Some of the key challenges include:

- **Large Volume of Data:** HPC systems generate terabytes of log data daily, making it difficult to process and analyze data in real time.
- **Sparse Failure Events:** Failures are relatively rare compared to normal system operations, leading to an imbalanced dataset that affects model training.
- **Complex Failure Patterns:** Failures may result from multiple contributing factors, requiring advanced data analysis techniques to identify correlations.
- **Interpretability of Predictions:** Machine learning models often operate as black boxes, making it difficult for administrators to understand why a specific failure prediction is made.
- **Real-Time Processing Requirements:** The failure prediction system must operate in real time to be effective, requiring efficient data pipelines and computing resources.

Addressing these challenges requires a combination of advanced log analysis, machine learning, and real-time monitoring to create a reliable failure prediction system.

4. Methods for HPC Failure Prediction

1. Data Collection and Preprocessing

The first step in developing an HPC failure prediction system is data collection. Various types of data sources contribute to failure prediction:

- **System Logs:** Kernel logs, syslogs, application logs, and security logs contain valuable information on system performance and errors.
- **Performance Metrics:** CPU utilization, memory usage, disk I/O, and network throughput help in assessing node health.
- **Hardware Sensor Data:** Temperature sensors, fan speed, and power consumption provide early warning signs of hardware degradation.
- **User Job Logs:** Information about job execution failures, resource allocation, and execution time contributes to identifying workload-induced failures.

Once the data is collected, preprocessing steps such as filtering, normalization, and feature extraction are performed. Removing redundant or irrelevant logs, handling missing values, and converting categorical data into numerical formats are essential steps for preparing the dataset for analysis.

2. Machine Learning Techniques for Failure Prediction

Machine learning techniques play a vital role in failure prediction. Several models can be used depending on the complexity of the data and prediction requirements:

- **Supervised Learning:** Techniques like Random Forest, Support Vector Machines (SVM), and Neural Networks are trained on labeled failure data to classify failure events.
- **Unsupervised Learning:** Clustering algorithms such as K-Means and DBSCAN detect anomalies by identifying deviations from normal system behavior.
- **Time Series Analysis:** Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) networks analyze time-dependent patterns in system logs to predict failures.
- **Anomaly Detection:** Autoencoders and Isolation Forests help detect outliers in monitoring metrics, signaling potential failures.

Feature engineering is crucial in improving model accuracy. Selecting the most relevant features, such as temperature spikes, memory leaks, and excessive disk errors, enhances the prediction performance.

3. Real-Time Monitoring and Prediction System

To make failure prediction effective, real-time monitoring and alerting mechanisms are necessary. A well-integrated system consists of:

- **Log Aggregation:** Tools like Elasticsearch, Logstash, and Kibana (ELK Stack) collect and visualize logs for better analysis.
- **Monitoring Systems:** Prometheus and Grafana provide real-time metric monitoring and alerting.
- **Automated Alerts:** Alerting mechanisms using email, SMS, or dashboards notify administrators about potential failures.
- **Predictive Maintenance:** Automated scripts trigger maintenance actions such as migrating workloads, restarting nodes, or isolating faulty hardware before failures occur.

MOTIVATION

HPC clusters generate **terabytes of log data and monitoring metrics** daily, which contain valuable insights into system behavior and failure patterns. However, analyzing this vast amount of data manually is impractical. With advancements in **machine learning (ML) and data analytics**, automated systems can be developed to detect anomalies and predict potential failures before they occur.

Why Failure Prediction is Essential?

HPC systems are widely used in critical domains such as weather forecasting, medical research, financial modeling, and large-scale simulations. Failures in these systems can lead to incomplete computations, data corruption, or even financial losses. As HPC workloads continue to grow in complexity, traditional reactive approaches to failure detection are no longer sufficient. Predictive maintenance is emerging as a necessity to ensure smooth and uninterrupted operations in these high-stakes environments.

Key Motivations for This Project

- **Minimize Downtime:**
 - Ensure continuous availability of HPC resources by proactively identifying failure risks.
 - Reduce unexpected system crashes that disrupt research and computations.
 - Improve job scheduling efficiency by predicting node failures in advance.
- **Enhance System Reliability:**
 - Implement intelligent monitoring to detect early warning signs of failures.
 - Reduce the frequency of unexpected crashes and ensure stable HPC performance.
 - Provide administrators with real-time insights into system health.
- **Improve Resource Utilization:**

- Optimize workload distribution based on predictive insights.
- Prevent resource wastage by avoiding node failures during critical computations.
- Ensure that computing power is allocated efficiently to meet performance demands.

- **Enable Data-Driven Decision Making:**

- Use real-time monitoring data to enhance administrative decisions.
- Implement an AI-driven approach for dynamic failure prediction and response.
- Assist system administrators in planning maintenance schedules based on predictive analytics.

- **Reduce Maintenance Costs:**

- Prevent costly hardware failures through predictive maintenance strategies.
- Reduce the need for emergency replacements and unplanned maintenance efforts.
- Extend the lifespan of HPC components by addressing potential failures before they occur.

Broader Implications

By addressing these challenges, the failure prediction system developed in this project will contribute to improving the overall efficiency of HPC environments. With reliable failure detection mechanisms, organizations can significantly reduce operational risks, increase computational efficiency, and ensure the smooth functioning of mission-critical applications. The integration of ML-driven insights into monitoring workflows will not only optimize resource management but also pave the way for more intelligent and autonomous HPC system administration.

PROBLEM STATEMENT

High-Performance Computing (HPC) environments are essential for executing complex computations and simulations across various scientific and industrial domains. However, these environments face significant challenges in detecting and mitigating node failures before they adversely impact system performance and overall productivity. Currently, most failure detection approaches are reactive, relying heavily on post-failure diagnostics. This means that issues are only addressed after they occur, resulting in unnecessary downtime, reduced efficiency, and potential data loss. The absence of automated failure prediction systems that can analyze historical log data alongside real-time metrics further exacerbates this problem, leaving HPC administrators with limited tools to preemptively address potential failures.

Key Challenges Addressed in This Project:

1. **Data Overload:** HPC clusters generate enormous amounts of log and monitoring data daily, encompassing everything from system performance metrics to error logs. The sheer volume of this data can be overwhelming, making it difficult for administrators to sift through and extract meaningful insights. Traditional log analysis techniques often fall short in this context, as they may not be equipped to handle the scale and complexity of the data generated. Advanced data processing techniques, including machine learning and big data analytics, are necessary to effectively manage and interpret this vast dataset.
2. **Reactive vs. Proactive Maintenance:** The prevailing maintenance strategies in HPC environments primarily focus on addressing failures after they occur. This reactive approach not only leads to unnecessary downtime but also increases operational costs and can hinder research progress. By developing a predictive failure model, this project aims to shift the maintenance paradigm towards a proactive approach, allowing for timely interventions that can prevent failures before they disrupt operations. This shift could significantly enhance system reliability and performance.

3. **Anomaly Detection:** Identifying what constitutes normal behavior for HPC nodes is inherently complex due to the dynamic nature of workloads, varying hardware conditions, and diverse job execution patterns. Anomalies can arise from a multitude of factors, including software bugs, hardware degradation, or unexpected workload spikes. Developing robust anomaly detection algorithms that can accurately discern between normal fluctuations and genuine issues is crucial for effective failure prediction. This requires a nuanced understanding of the operational context and the ability to adapt to changing conditions.
4. **Real-Time Monitoring and Alerts:** There is a pressing need to integrate failure prediction capabilities with existing real-time monitoring tools. This integration would enable the generation of actionable alerts that can inform administrators of potential issues before a node fails. By providing timely notifications, HPC environments can implement corrective actions proactively, thereby minimizing downtime and maintaining system performance. The challenge lies in ensuring that these alerts are not only timely but also relevant and actionable.
5. **Model Selection and Optimization:** The selection of appropriate machine learning techniques for anomaly detection and predictive modeling in HPC environments is a critical challenge. With a plethora of algorithms available, determining which models are best suited for the unique characteristics of HPC data is essential. Additionally, optimizing these models for performance and accuracy in real-world scenarios is necessary to ensure their effectiveness. This involves not only selecting the right algorithms but also fine-tuning their parameters and validating their performance against historical data.

OBJECTIVES

The primary goal of this project is to develop an effective HPC Node/Cluster Failure Prediction System that leverages logs and monitoring metrics to proactively detect potential failures. The overarching goal is to enhance the reliability, efficiency, and maintainability of High-Performance Computing (HPC) clusters, ultimately reducing operational costs and ensuring uninterrupted service availability. The specific objectives of the project are outlined below:

1. Proactive Maintenance

One of the primary objectives is to identify potential failures before they occur, thereby minimizing downtime. By implementing a predictive maintenance strategy, the project seeks to shift the focus from reactive responses to proactive interventions. This approach will allow administrators to address issues before they escalate into significant problems, ensuring continuous operation and optimal performance of HPC resources.

2. Log Collection and Analysis

The project will involve comprehensive log analysis, utilizing historical log data to detect patterns that precede node or cluster failures. This includes:

- **Log Collection:** Systematically gathering logs from various sources, including system logs, job scheduler logs (e.g., SLURM, PBS), and hardware metrics (CPU usage, memory consumption, disk I/O, network traffic).
- **Pattern Detection:** Analyzing the collected logs to identify recurring patterns or anomalies that may indicate impending failures. This step is crucial for building a robust predictive model.

3. Real-Time Monitoring

Incorporating real-time monitoring is essential for the success of the failure prediction system. The project will focus on:

- **Metrics Integration:** Monitoring key performance indicators such as CPU usage, memory consumption, network traffic, and disk health. These metrics will provide valuable insights into the operational state of the HPC cluster.
- **Anomaly Detection:** Implementing real-time anomaly detection mechanisms to identify deviations from normal behavior, which could signal potential failures.

4. Machine Learning-Based Failure Prediction

The core of the project lies in developing machine learning models that can classify or regress failure probabilities based on both historical and real-time data. This objective includes:

- **Model Development:** Training various predictive models, including anomaly detection algorithms (e.g., Isolation Forest, DBSCAN) and predictive models (e.g., XGBoost, LSTMs, RNNs).
- **Performance Evaluation:** Evaluating the models using metrics such as Precision, Recall, F1 Score, and AUC-ROC to determine their effectiveness in predicting failures.

5. Integration with Monitoring Systems

To ensure that the predictive model is actionable, the project will focus on integrating it with existing monitoring frameworks. This includes:

- **Deployment within Monitoring Tools:** Integrating the failure prediction model with monitoring solutions like Prometheus and Grafana to facilitate real-time data analysis and visualization.
- **Automated Alerting System:** Implementing an automated alerting system that notifies administrators of potential failures based on the predictions made by the model. This will enable timely interventions and proactive maintenance.

6. Scalable Deployment

To ensure that the failure prediction system can be effectively utilized in various HPC environments, the project will focus on scalable deployment. This involves:

- **Containerization:** Utilizing containerization technologies such as Docker to package the predictive model and its dependencies, ensuring that it can be easily deployed across different environments.
- **Orchestration with Kubernetes:** Implementing Kubernetes for orchestration, allowing for the management of containerized applications at scale.

7. Comparison of Predictive Models

An important objective of the project is to benchmark the performance of different machine learning models to determine the most effective approach for failure prediction. This involves:

- **Model Evaluation:** Conducting a systematic comparison of various predictive models, including traditional machine learning algorithms and advanced deep learning techniques.
- **Hyperparameter Tuning:** Implementing techniques to optimize the hyperparameters of each model, ensuring that they perform at their best.

8. Implementation of a Decision Support System

To enhance the usability of the predictive failure system, the project will develop a Decision Support System (DSS). This objective includes:

- **Dashboard Development:** Creating an intuitive dashboard or visualization tool that provides administrators with real-time insights into the health of the HPC cluster.
- **Insights and Recommendations:** Incorporating analytical tools that provide actionable insights and recommendations for proactive maintenance based on the predictive model's outputs.

9. Performance Testing and Optimization

To ensure that the failure prediction system operates effectively in real-world conditions, extensive performance testing and optimization will be conducted. This objective involves:

- **System Performance Testing:** Conducting rigorous testing to measure the system's performance under various load conditions.

- **Model Optimization:** Continuously refining the predictive models based on performance testing results.

10. Deployment in a Real HPC Environment (Continued)

- **Model Refinement:** Gathering feedback from the pilot deployment to refine the predictive models and the overall system based on observed results. This iterative process will help ensure that the system is tailored to the specific needs and characteristics of the HPC environment.
- **Documentation and Training:** Providing comprehensive documentation and training for HPC administrators on how to use the system effectively. This will include guidelines on interpreting alerts, utilizing the Decision Support System, and performing maintenance based on predictive insights.

By deploying the system in a real HPC environment, the project aims to demonstrate its practical applicability and effectiveness in enhancing the reliability and efficiency of HPC operations.

PLANNING

Introduction

Planning is a pivotal phase in the software development process, setting the foundation for subsequent stages. It involves strategizing the approach, determining tasks, scheduling, and resource allocation to ensure the effective completion of the project.

Planning Objectives

The primary objective of planning is to identify and schedule tasks necessary for project completion. A robust plan should be flexible enough to accommodate unforeseen events while considering economic, political, and human factors. While a detailed requirements document is not mandatory for planning, awareness of crucial requirements is essential.

Major Issues Addressed in Project Plan:-

1. **Cost Estimation:** Estimating the financial resources required for project execution.
2. **Schedule and Milestones:** Establishing a timeline with key milestones to track progress.
3. **Personnel Plan:** Allocating human resources effectively based on skill sets and availability.
4. **Software Quality Assurance Plans:** Outlining strategies to ensure software quality throughout the project.
5. **Configuration Management Plans:** Defining procedures for managing changes to the software configuration.
6. **Project Monitoring Plans:** Establishing mechanisms for monitoring project progress and addressing deviations.
7. **Risk Management:** Identifying potential risks and developing strategies to mitigate them.

Planning Activities

Planning activities include establishing the need for the system, conducting sensitivity assessments, performing initial risk assessments, and reviewing solicitation documents. These activities ensure that the project plan addresses all critical aspects and aligns with project objectives.

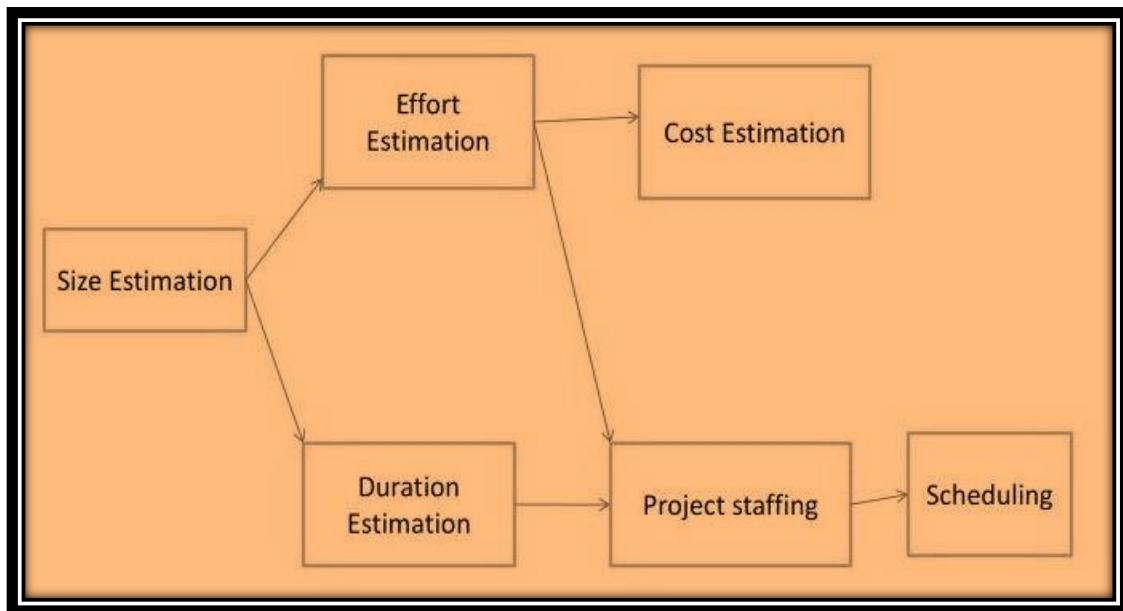


Fig.6.1–Planning activities

Waterfall Model

The Waterfall Model outlines the sequential execution of project phases:

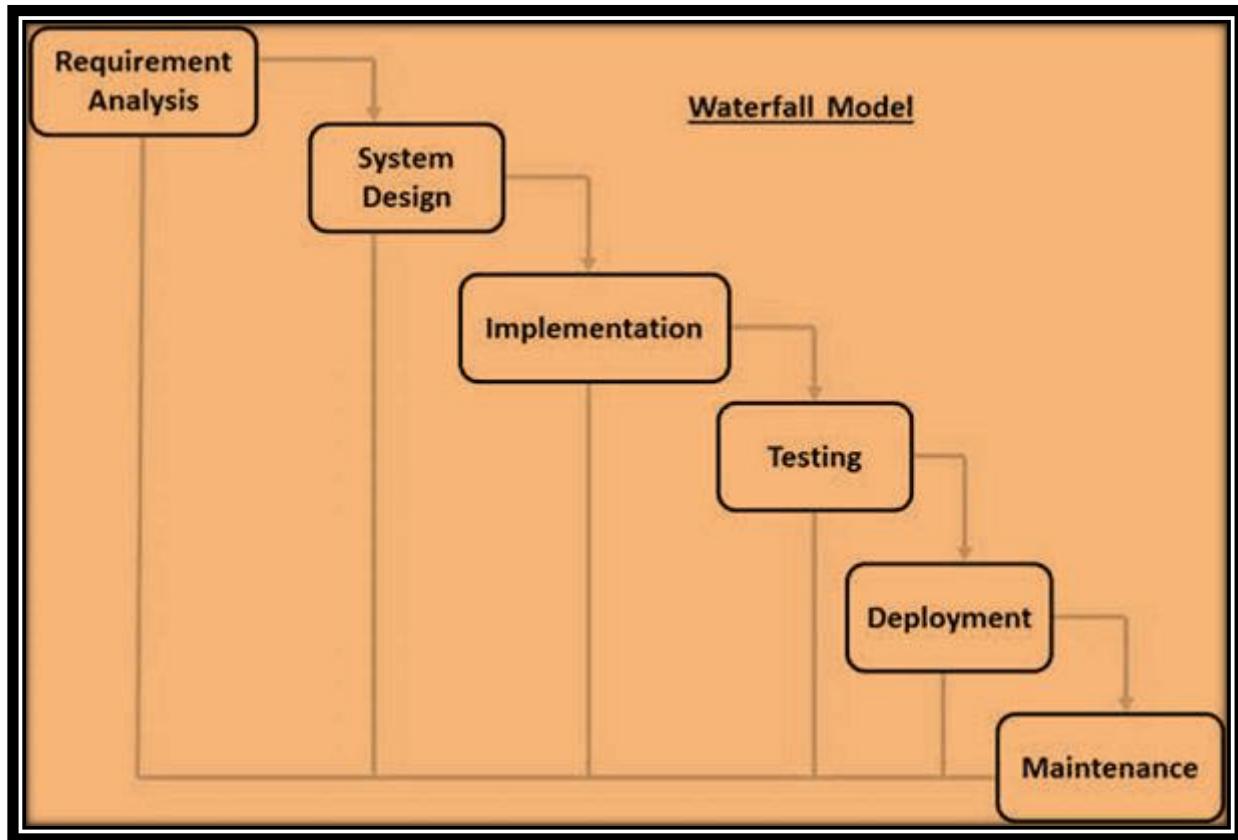


Fig.6.2–Waterfall model

- Gathering and analysis of needs – During this stage, all potential system requirements are gathered and recorded in a requirements specification document.
- System Design – In this step, the required specifications from the first phase are examined, and a system design is created. This system design determines the overall system architecture as well as the hardware and system requirements.
- Implementation The system is initially developed as discrete programmes known as units, which are then combined in the next phase, with input from the system design. Unit testing is the process of developing and evaluating each unit for functionality.

- Integration and Testing – All the units developed in the implementation phase are integrated into a system after testing of each unit. Post integration the entire system is tested for any faults and failures.
- Deployment of system – Once the functional and non-functional testing is done; the product is deployed in the customer environment or released into the market.
- Maintenance – There are some issues which come up in the client environment. To fix those issues, patches are released. Also to enhance the product some better versions are released. Maintenance is done to deliver these changes in the customer environment.

Planning serves as the roadmap for successful project execution, ensuring that all aspects of the project are meticulously addressed and executed in a timely manner. It sets the stage for subsequent phases, guiding the development team towards project completion.

METHODOLOGY

Methodology Overview

The proposed failure prediction system follows a structured workflow designed to systematically address the challenges of detecting and mitigating node failures in High-Performance Computing (HPC) environments. This methodology encompasses several key steps, each critical to the successful development and implementation of the predictive system. Below is a detailed elaboration of each step in the workflow.

Step 1: Data Collection

The first step in the methodology involves comprehensive data collection from various sources within the HPC environment. This step is crucial for building a robust dataset that will serve as the foundation for subsequent analysis and modeling.

- **System Logs:** Collect logs from multiple nodes, including:
 - **Job Scheduler Logs:** These logs (e.g., from SLURM or PBS) provide insights into job execution, resource allocation, and scheduling delays. They are essential for understanding how jobs interact with the system and can reveal patterns leading to failures.
 - **System Event Logs:** These logs capture system-level events, including warnings, errors, and other significant occurrences that may indicate underlying issues.
 - **Application Logs:** Logs generated by applications running on the HPC cluster can provide context on how specific workloads may contribute to failures.
- **Real-Time Hardware Metrics:** Monitor critical hardware metrics using tools like Prometheus and Node Exporter. Key metrics to be collected include:
 - **CPU/GPU Utilization:** Monitoring the usage levels of processing units helps identify resource bottlenecks.
 - **Memory Usage:** Tracking memory consumption is vital for detecting potential memory leaks or overloads.
 - **Disk I/O:** Monitoring disk input/output operations can reveal performance issues related to storage.
 - **Network Traffic:** Analyzing network traffic helps identify potential bottlenecks or failures in communication between nodes.
- **Failure Event Identification:** Identify and annotate failure events within the collected logs.

This involves defining what constitutes a failure in the context of the HPC environment and tagging relevant log entries accordingly. This annotated data will be critical for training predictive models.

Step 2: Data Preprocessing & Feature Engineering

Once the data is collected, the next step is to preprocess it and engineer features that will enhance the predictive capabilities of the models.

- **Log Parsing:** Use tools such as Fluentd, Logstash, or custom Python scripts to parse the collected logs. This process involves:
 - Structuring unstructured log data into a format suitable for analysis.
 - Filtering out irrelevant information and focusing on entries that are pertinent to failure prediction.
- **Feature Extraction:** Extract key features from the logs and metrics that may indicate potential failures. Important features to consider include:
 - **Error Message Frequency:** The frequency of specific error messages can indicate underlying issues.
 - **CPU Temperature Trends:** Monitoring temperature trends can help identify overheating issues that may lead to hardware failures.
 - **I/O Failure Counts:** Tracking the number of I/O failures can provide insights into potential disk-related issues.
- **Data Cleaning:** Handle missing or redundant data to ensure model accuracy. This may involve:
 - Imputing missing values using statistical methods or removing incomplete records.
 - Eliminating duplicate entries to prevent skewing the analysis.

Step 3: Exploratory Data Analysis (EDA)

Exploratory Data Analysis is a critical step for understanding the data and identifying patterns that may inform the predictive modeling process.

- **Pattern Identification:** Use visualization tools such as Matplotlib, Seaborn, Kibana, and Grafana to identify patterns and correlations in failure events. This includes:
 - Visualizing trends over time to see how metrics change leading up to failures.
 - Analyzing the distribution of features to understand their behavior under normal and failure conditions.
- **Anomaly Detection:** Detect anomalies in log sequences that could indicate potential failures.

This involves:

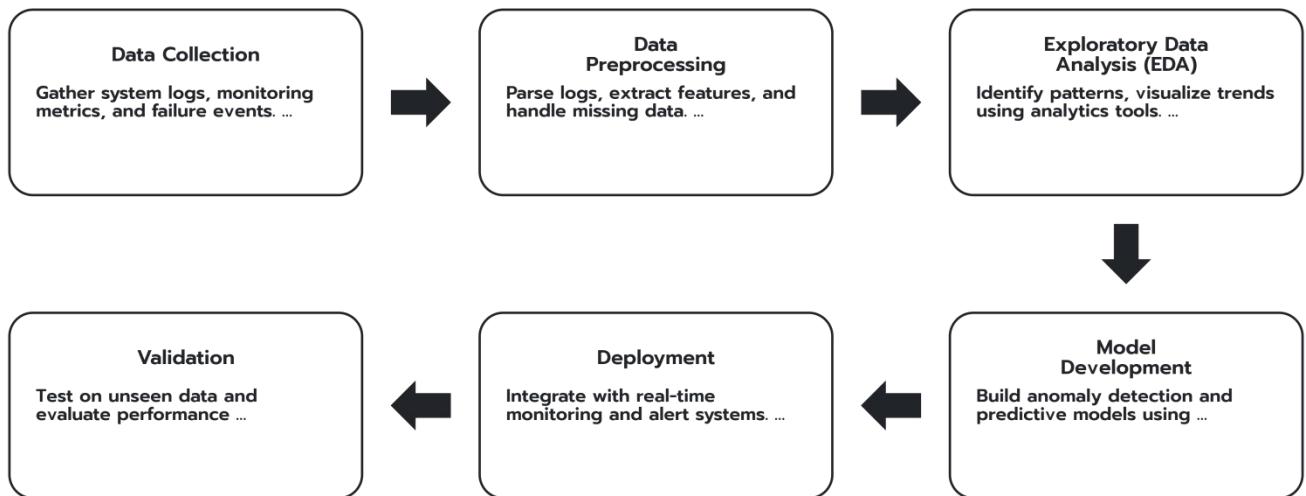
- Applying statistical methods to identify outliers in the data.
- Using visualization techniques to highlight unusual patterns that may warrant further investigation.

Step 4: Machine Learning-Based Prediction

The core of the methodology involves developing machine learning models to predict failures based on the processed data.

- **Anomaly Detection Models:** Train models such as Isolation Forest and DBSCAN to identify unusual behaviors in the data. These models will help flag potential anomalies that could indicate impending failures.
- **Supervised Learning Models:** Develop predictive models using algorithms such as XGBoost, Random Forest, and Long Short-Term Memory (LSTM) networks. This step includes:
 - Training the models on the annotated dataset to classify or regress failure probabilities based on historical and real-time data.
 - Evaluating model performance using metrics such as Precision, Recall, F1 Score, and AUC-ROC to ensure the models are effective in predicting failures.

Workflow

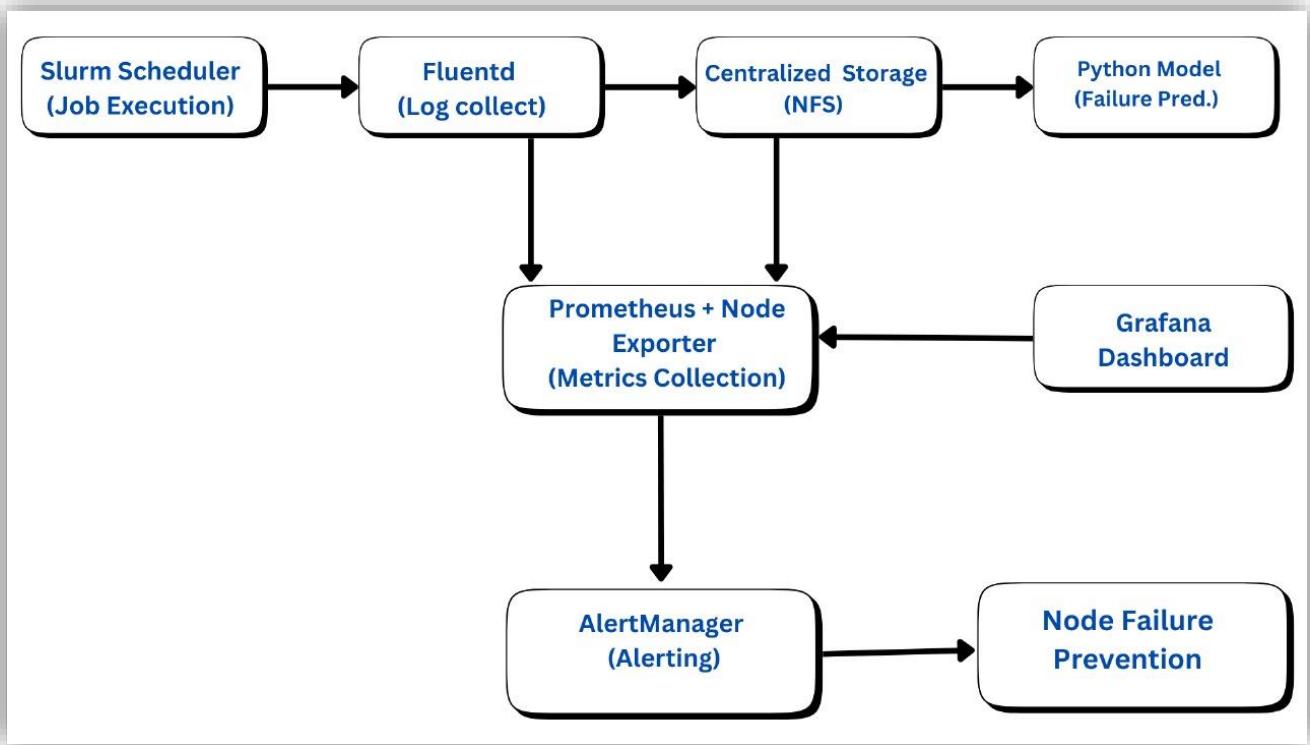


1. **Data Collection:** Gather system logs, monitoring metrics, and failure events from various sources.
2. **Data Preprocessing:** Parse logs, extract relevant features, and handle missing or redundant data to prepare the dataset for analysis.
3. **Exploratory Data Analysis (EDA):** Identify patterns, visualize trends, and detect anomalies using analytics tools to gain insights into the data.
4. **Model Development:** Build and train anomaly detection and predictive models using machine learning and deep learning algorithms to forecast potential failures.
5. **Deployment:** Integrate the predictive models with real-time monitoring and alert systems to provide timely notifications to administrators.
6. **Validation:** Test the models on unseen data to evaluate their performance and ensure they meet the required accuracy and reliability standards.

This structured methodology not only facilitates a comprehensive understanding of the data but also ensures that the predictive system is robust, scalable, and capable of adapting to the dynamic nature of HPC environments. By following these steps, the project aims to create a reliable failure prediction system that enhances the operational efficiency and reliability of HPC clusters.

IMPLEMENTATION STEPS

The HPC Node/Cluster Failure Prediction System is designed to proactively identify potential failures in high-performance computing environments by leveraging logs and monitoring metrics. This system aims to enhance the reliability and efficiency of HPC clusters, ultimately reducing downtime and operational costs.



1. System Requirements

Hardware Requirements

- **Master Node:** Minimum 8 CPU cores, 32 GB RAM, and 500 GB SSD.
- **Compute Nodes:** Minimum 4 CPU cores, 16 GB RAM, and 250 GB SSD per node.
- **Networking:** High-speed interconnects (e.g., InfiniBand or Gigabit Ethernet).
- **Storage:** Shared storage solution (e.g., NFS or Lustre).

Software Requirements

- **Operating System:** Linux distribution (e.g., CentOS, Ubuntu Server).
- **Cluster Management Software:** OpenHPC or Rocks Cluster.
- **Job Scheduler:** SLURM, PBS, or Torque.
- **Monitoring Tools:** Prometheus, Node Exporter, Grafana.
- **Log Analysis Tools:** Fluentd, Logstash, Elasticsearch.
- **Machine Learning Libraries:** Python (Pandas, NumPy, Scikit-learn, TensorFlow/PyTorch).

2. Installation and Setup

Setting Up the HPC Cluster

1. **Hardware Setup:** Install the master node and compute nodes as per the hardware requirements.
2. **Networking:** Configure high-speed networking between nodes.
3. **Storage Configuration:** Set up shared storage using NFS or Lustre.

Installing Required Software

1. **Operating System:** Install the chosen Linux distribution on all nodes.
2. **Cluster Management:** Install OpenHPC or Rocks Cluster for managing the HPC environment.
3. **Job Scheduler:** Install SLURM, PBS, or Torque for job scheduling.

Configuring Log Collection and Monitoring

1. **Install Prometheus** on the master node for metric collection.
2. **Install Node Exporter** on all compute nodes to expose hardware metrics.
3. **Set Up Fluentd or Logstash** for centralized log collection.

3. Data Collection

Collecting System Logs

1. **Centralize Log Files:** Use rsyslog or Fluentd to collect logs from:

- Job scheduler (e.g., SLURM logs).
- System logs (/var/log/syslog, /var/log/messages).
- Application logs.

Monitoring Hardware Metrics

1. **Configure Prometheus** to scrape metrics from Node Exporter.
2. **Monitor Key Metrics:** CPU/GPU utilization, memory usage, disk health, and network traffic.

Annotating Failure Events

1. **Define Failure Events:** Identify what constitutes a failure (e.g., hardware failure, job crashes).
2. **Annotate Logs:** Tag relevant log entries with failure events for future analysis.

4. Data Preprocessing

Log Parsing

1. **Use Fluentd or Logstash** to parse collected logs.
2. **Extract Useful Information:** Filter out irrelevant data and focus on entries pertinent to failure prediction.

Feature Engineering

1. **Count Error Messages:** Track occurrences of specific error messages (e.g., I/O errors, kernel panics).
2. **Compute Resource Trends:** Analyze trends in CPU/GPU temperature, memory usage, and disk I/O.
3. **Identify Anomalous Patterns:** Use statistical methods to detect outliers in the data.

Handling Missing Data

1. **Identify Missing Entries:** Use data profiling techniques to find gaps in the dataset.
2. **Impute or Remove:** Decide whether to fill in missing values using techniques like mean imputation or to remove incomplete records.

5. Exploratory Data Analysis (EDA)

Tools for EDA

- **Matplotlib:** For creating static, animated, and interactive visualizations in Python.
- **Seaborn:** For statistical data visualization based on Matplotlib.
- **Grafana:** For real-time monitoring and visualization of metrics.

Identifying Patterns and Correlations

1. **Visualize Data:** Create plots to visualize trends in log frequency and resource usage.
2. **Correlation Analysis:** Use heatmaps to identify correlations between different metrics and failure events.

6. Model Development

Training Anomaly Detection Models

1. **Select Algorithms:** Choose from Isolation Forest, DBSCAN, or PCA for unsupervised anomaly detection.

2. **Train Models:** Use historical logs and metrics to train the models.

Developing Predictive Models

1. **Select Algorithms:** Use Random Forest, XGBoost, or LSTM for supervised learning.
2. **Train and Validate:** Split the dataset into training and testing sets, and validate model performance.

Model Evaluation

1. **Performance Metrics:** Evaluate models using Precision, Recall, F1 Score, and AUC-ROC.
2. **Hyperparameter Tuning:** Optimize model parameters for better performance.

7. Deployment

Integrating with Real-Time Monitoring Tools

1. **Deploy the Predictive Model:** Use TensorFlow Serving or ONNX for model deployment.
2. **Integrate with Grafana:** Set up Grafana to visualize predictions and metrics.

Setting Up Alert Systems

1. **Configure AlertManager:** Set up alerts based on failure probability thresholds.
2. **Notification Channels:** Integrate with email, Slack, or PagerDuty for alert notifications.

Containerization and Scalability

1. **Use Docker:** Containerize the application for easy deployment.
2. **Kubernetes:** Deploy the application on Kubernetes for scalability and management.

8. Validation and Testing

Testing the Model on Unseen Data

1. **Simulate Node Failures:** Conduct tests by simulating failures to evaluate the model's predictions.
2. **Analyze Results:** Review the model's performance on unseen data to ensure reliability.

Performance Metrics

1. **Evaluate Metrics:** Use Precision, Recall, and F1 Score to assess model effectiveness.
2. **Iterate:** Refine the model based on testing outcomes.

9. Using the Dashboard

Accessing the Dashboard

1. **Open Grafana:** Navigate to the Grafana URL in a web browser.
2. **Log In:** Use the credentials set during installation.

Interpreting Alerts and Predictions

1. **View Real-Time Metrics:** Monitor the dashboard for live updates on node status.
2. **Understand Alerts:** Review alerts to identify potential failures and take preemptive actions.

10. Troubleshooting

Common Issues and Solutions

1. **Log Collection Failures:** Ensure that Fluentd or Logstash is running and properly configured.
2. **Model Performance Issues:** Revisit feature engineering and model selection if predictions are inaccurate.
3. **Alerting Problems:** Check AlertManager configuration and notification channels for issues.

SCREENSHOTS

Step-by-Step Implementation Screenshots

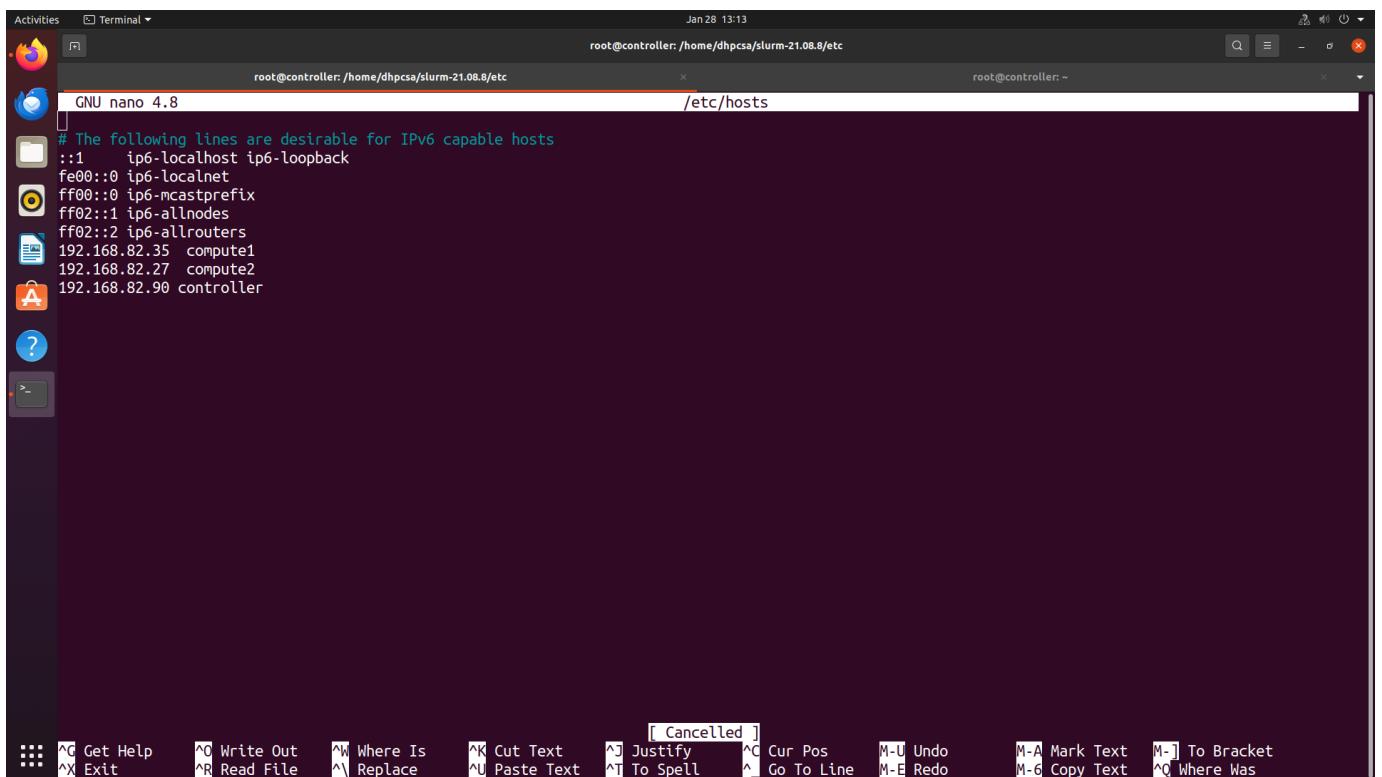
The following sections provide a detailed overview of the implementation process for the HPC Node/Cluster Failure Prediction System. Each step includes a brief description and a corresponding screenshot to illustrate the process.

1. SLURM Installation

Description: SLURM (Simple Linux Utility for Resource Management) is a job scheduler used to manage compute resources in HPC environments. This step involves installing and configuring SLURM on the master and compute nodes.

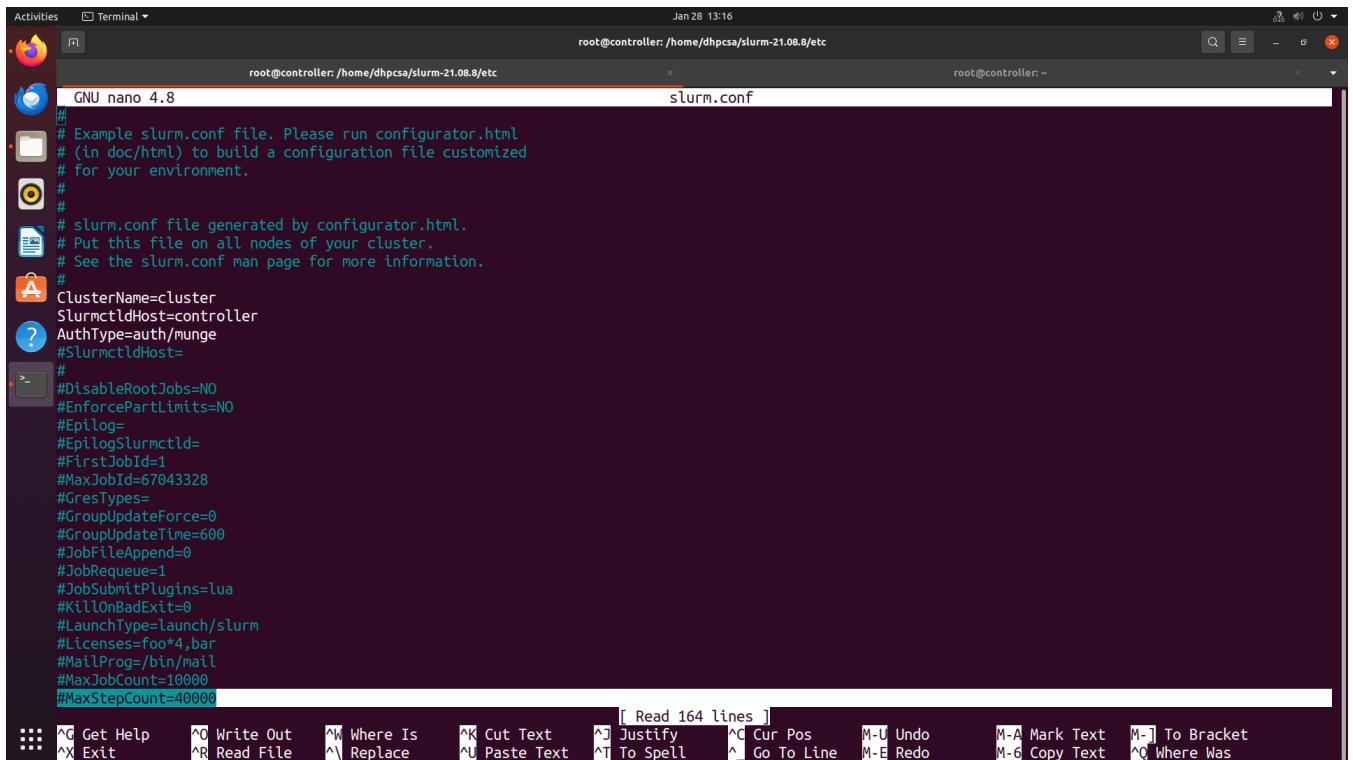
Screenshot:

etc/hosts file



```
# The following lines are desirable for IPv6 capable hosts
::1      ip6-localhost ip6-loopback
fe00::0  ip6-mcastprefix
ff02::1  ip6-allnodes
ff02::2  ip6-allrouters
192.168.82.35  compute1
192.168.82.27  compute2
192.168.82.90  controller
```

SLURM.CONF FILE



```

# Example slurm.conf file. Please run configurator.html
# (in doc/html) to build a configuration file customized
# for your environment.
#
# slurm.conf file generated by configurator.html.
# Put this file on all nodes of your cluster.
# See the slurm.conf man page for more information.
#
# ClusterName=cluster
SlurmctldHost=controller
AuthType=auth/munge
#SlurmctldHost=
#
#DisableRootJobs=NO
#EnforcePartLimits=NO
#Epilog=
#EpilogSlurmctld=
#FirstJobId=1
#MaxJobId=67043328
#GresTypes=
#GroupUpdateForce=0
#GroupUpdateTime=600
#JobfileAppend=0
#JobRequeue=1
#JobSubmitPlugins=lua
#KillOnBadExit=0
#LaunchType=launch/slurm
#Licenses=foo*4,bar
#MailProg=/bin/mail
#MaxJobCount=10000
#MaxStepCount=40000

```

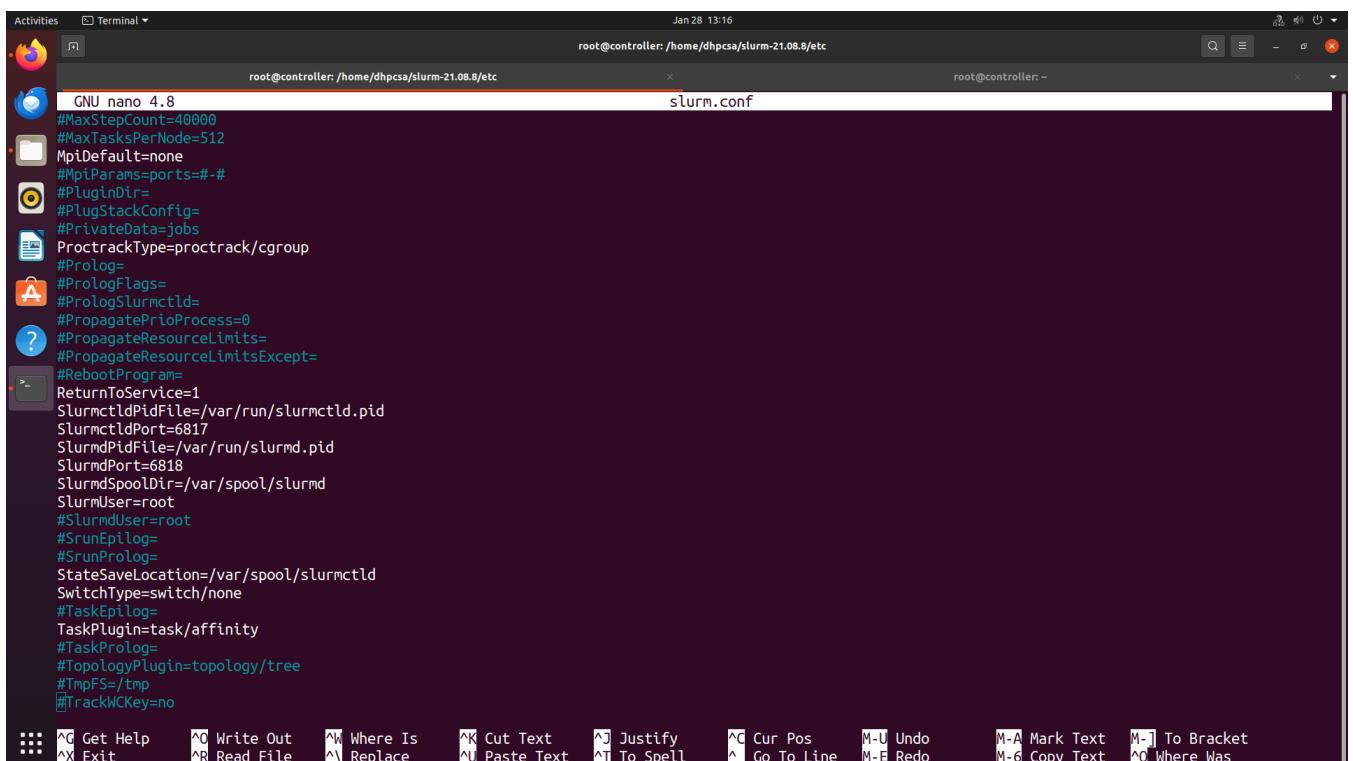
[Read 164 lines]

GNU nano 4.8

slurm.conf

Activities Terminal Jan 28 13:16 root@controller: /home/dhpcsa/slurm-21.08.8/etc root@controller: ~

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify ^C Cur Pos M-U Undo M-A Mark Text M-[To Bracket
 ^X Exit ^R Read File ^A Replace ^U Paste Text ^T To Spell ^I Go To Line M-E Redo M-6 Copy Text ^Q Where Was



```

#MaxStepCount=40000
#MaxTasksPerNode=512
MpDefault=None
#MpParams=ports=-#
#PluginDir=
#PlugStackConfig=
#PrivateData=jobs
#ProctrackType=proctrack/cgroup
#Prolog=
#PrologFlags=
#PrologSlurmctld=
#PropagatePrioProcess=0
#PropagateResourceLimits=
#PropagateResourceLimitsExcept=
#RebootProgram=
ReturnToService=1
SlurmctldPidfile=/var/run/slurmctld.pid
SlurmctldPort=6817
SlurmdPidfile=/var/run/slurmd.pid
SlurmdPort=6818
SlurmdSpoolDir=/var/spool/slurmd
SlurmUser=root
#SlurmUser=root
#SrunEpilog=
#SrunProlog=
StateSaveLocation=/var/spool/slurmctld
SwitchType=switch/none
#TaskEpilog=
TaskPlugin=task/affinity
#TaskProlog=
#TopologyPlugin=topology/tree
#TmpFS=/tmp
#TrackWCKey=no

```

[Read 164 lines]

GNU nano 4.8

slurm.conf

Activities Terminal Jan 28 13:16 root@controller: /home/dhpcsa/slurm-21.08.8/etc root@controller: ~

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify ^C Cur Pos M-U Undo M-A Mark Text M-[To Bracket
 ^X Exit ^R Read File ^A Replace ^U Paste Text ^T To Spell ^I Go To Line M-E Redo M-6 Copy Text ^Q Where Was

HPC Node/Cluster Failure Prediction With Prevention



```
GNU nano 4.8
#TrackKey=no
#TreeWidth=
#UnkillableStepProgram=
#UsePAM=0
#
#
# TIMERS
#BatchStartTimeout=10
#CompleteWait=0
#EpilogMsgTime=2000
#GetEnvTimeout=2
#HealthCheckInterval=0
#HealthCheckProgram=
InactiveLimit=0
KillWait=30
#MessageTimeout=10
#ResvOverRun=0
MinJobAge=300
#OverTimeLimit=0
SlurmctldTimeout=120
SlurmdTimeout=300
#UnkillableStepTimeout=60
#VSizeFactor=0
Waittime=0
#
#
# SCHEDULING
#DefMemPerCPU=0
#MaxMemPerCPU=0
#SchedulerTimeSlice=30
SchedulerType=sched/backfill
SelectType=select/cons_tres
SelectTypeParameters=CR_CPU_Memory
```

The screenshot shows a Linux desktop environment with a terminal window open in the background. The terminal window title is "root@controller: /home/dhpcsa/slurm-21.08.8/etc". The terminal content shows the "slurm.conf" file being edited in the "GNU nano 4.8" editor. The configuration file contains various parameters for Slurm, such as job priority, accounting storage, and job completion handling.

```
GNU nano 4.8
SelectTypeParameters=CR_CPU_Memory
PreemptType=preempt/qos
PreemptMode=suspend,gang
#
#
# JOB PRIORITY
#PriorityFlags=
PriorityType=priority/multifactor
#PriorityDecayHalfLife=
#PriorityCalcPeriod=
#PriorityFavorSmall=
#PriorityMaxAge=
#PriorityUsageResetPeriod=
#PriorityWeightAge=
#PriorityWeightFairshare=
#PriorityWeightJobSize=
#PriorityWeightPartition=
#PriorityWeightQOS=
#
#
# LOGGING AND ACCOUNTING
AccountingStorageEnforce=limits,qos
#AccountingStorageHost=
#AccountingStoragePass=
#AccountingStoragePort=
AccountingStorageType=accounting_storage/slurmdbd
#AccountingStorageUser=
#AccountingStorageFlags=
#JobCompHost=
#JobCompLoc=
#JobCompPass=
#JobCompPort=
#JobCompType=jobcomp/none
```

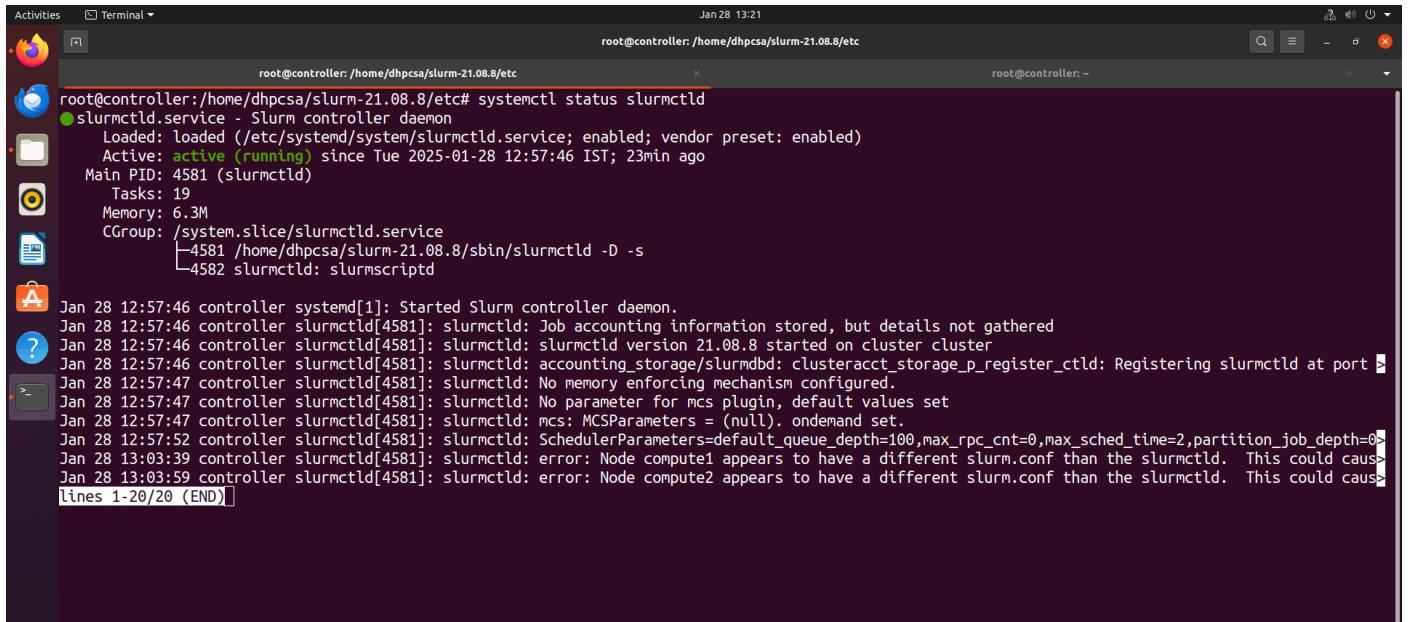
```
Activities Terminal Jan 28 13:17
root@controller:/home/dhpcsa/slurm-21.08.8/etc
root@controller: /home/dhpcsa/slurm-21.08.8/etc
root@controller: ~

GNU nano 4.8
JobCompType=jobcomp/none
#JobCompUser=
#JobContainerType=job_container/none
JobAccGatherFrequency=30
#JobAccGatherType=jobacct_gather/none
SlurmctldDebug=info
SlurmctldLogFile=/var/log/slurmctld.log
SlurmdDebug=info
SlurmdLogFile=/var/log/slurmd.log
#SlurmSchedLogFile=
#SlurmSchedLogLevel=
#DebugFlags=
#
#
# POWER SAVE SUPPORT FOR IDLE NODES (optional)
#SuspendProgram=
#ResumeProgram=
#SuspendTimeout=
#ResumeTimeout=
#ResumeRate=
#SuspendExcNodes=
#SuspendExcParts=
#SuspendRate=
#SuspendTime=
#
#
# COMPUTE NODES
NodeName=controller CPUs=4 State=UNKNOWN
NodeName=compute[1-2] CPUs=4 State=UNKNOWN

PartitionName=newpartition Nodes=ALL Default=YES MaxTime=02:00:00 State=UP
Mailprog=/usr/bin/mail
```

SCRIPT TO COPY SLURM.CONF FILE TO ALL COMPUTE NODES

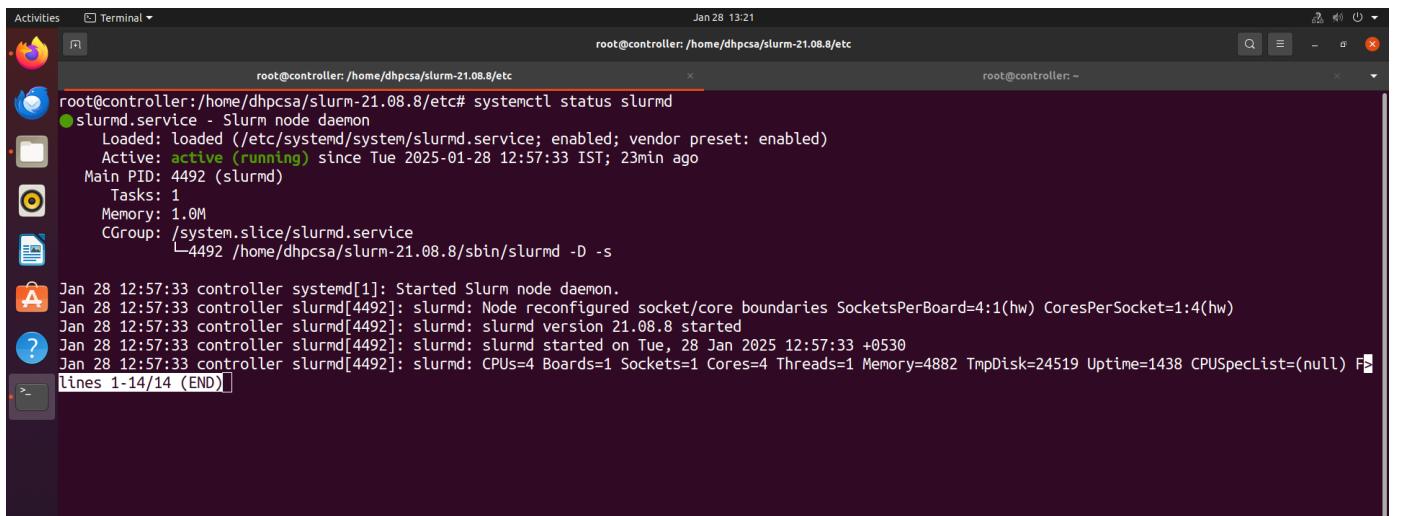
RUNNING STATUS OF SLURMCTLD DAEMON



```
root@controller:/home/dhpcsa/slurm-21.08.8/etc# systemctl status slurmctld
● slurmctld.service - Slurm controller daemon
   Loaded: loaded (/etc/systemd/system/slurmctld.service; enabled; vendor preset: enabled)
   Active: active (running) since Tue 2025-01-28 12:57:46 IST; 23min ago
     Main PID: 4581 (slurmctld)
        Tasks: 19
       Memory: 6.3M
      CGroup: /system.slice/slurmctld.service
              └─4581 /home/dhpcsa/slurm-21.08.8/sbin/slurmctld -D -s
              ├─4582 slurmctld: slurmscriptd

A Jan 28 12:57:46 controller systemd[1]: Started Slurm controller daemon.
Jan 28 12:57:46 controller slurmctld[4581]: slurmctld: Job accounting information stored, but details not gathered
? Jan 28 12:57:46 controller slurmctld[4581]: slurmctld: slurmctld version 21.08.8 started on cluster cluster
Jan 28 12:57:46 controller slurmctld[4581]: slurmctld: accounting_storage/slurmdbd: clusteracct_storage_p_register_ctld: Registering slurmctld at port 5000
Jan 28 12:57:47 controller slurmctld[4581]: slurmctld: No memory enforcing mechanism configured.
Jan 28 12:57:47 controller slurmctld[4581]: slurmctld: No parameter for mcs plugin, default values set
Jan 28 12:57:47 controller slurmctld[4581]: slurmctld: mcs: MCSParameters = (null). ondemand set.
Jan 28 12:57:52 controller slurmctld[4581]: slurmctld: SchedulerParameters=default_queue_depth=100,max_rpc_cnt=0,max_sched_time=2,partition_job_depth=0
Jan 28 13:03:39 controller slurmctld[4581]: slurmctld: error: Node compute1 appears to have a different slurm.conf than the slurmctld. This could cause problems.
Jan 28 13:03:59 controller slurmctld[4581]: slurmctld: error: Node compute2 appears to have a different slurm.conf than the slurmctld. This could cause problems.
[lines 1-20/20 (END)]
```

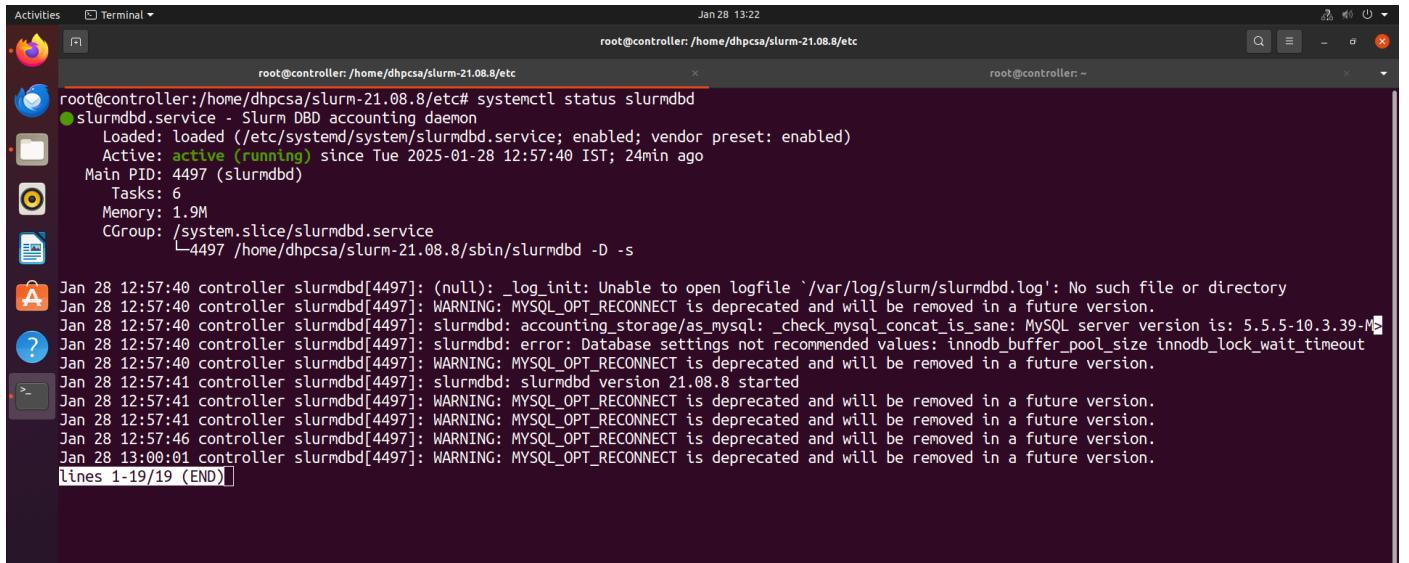
RUNNING STATUS OF SLURMD DAEMON



```
root@controller:/home/dhpcsa/slurm-21.08.8/etc# systemctl status slurmd
● slurmd.service - Slurm node daemon
   Loaded: loaded (/etc/systemd/system/slurmd.service; enabled; vendor preset: enabled)
   Active: active (running) since Tue 2025-01-28 12:57:33 IST; 23min ago
     Main PID: 4492 (slurmd)
        Tasks: 1
       Memory: 1.0M
      CGroup: /system.slice/slurmd.service
              └─4492 /home/dhpcsa/slurm-21.08.8/sbin/slurmd -D -s

A Jan 28 12:57:33 controller systemd[1]: Started Slurm node daemon.
Jan 28 12:57:33 controller slurmd[4492]: slurmd: Node reconfigured socket/core boundaries SocketsPerBoard=4:1(hw) CoresPerSocket=1:4(hw)
Jan 28 12:57:33 controller slurmd[4492]: slurmd: slurmd version 21.08.8 started
? Jan 28 12:57:33 controller slurmd[4492]: slurmd: slurmd started on Tue, 28 Jan 2025 12:57:33 +0530
Jan 28 12:57:33 controller slurmd[4492]: slurmd: CPUs=4 Boards=1 Sockets=4 Cores=4 Threads=1 Memory=4882 TmpDisk=24519 Uptime=1438 CPUSpecList=(null) F
[lines 1-14/14 (END)]
```

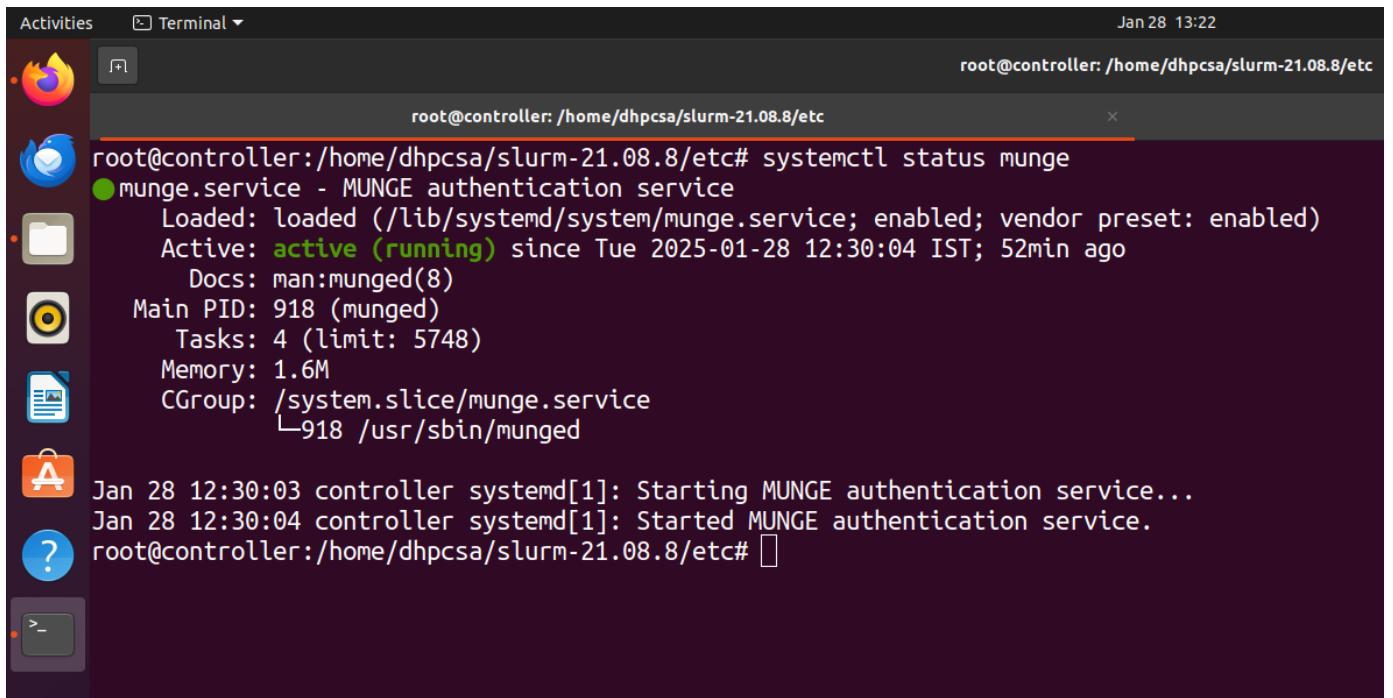
RUNNING STATUS OF SLURMDBD DAEMON



```
root@controller:/home/dhpcsa/slurm-21.08.8/etc# systemctl status slurmdbd
● slurmdbd.service - Slurm DBD accounting daemon
  Loaded: loaded (/etc/systemd/system/slurmdbd.service; enabled; vendor preset: enabled)
  Active: active (running) since Tue 2025-01-28 12:57:40 IST; 24min ago
    Main PID: 4497 (slurmdbd)
      Tasks: 6
     Memory: 1.9M
        CPU: 0.000 CPU(s) since start
       CGroup: /system.slice/slurmdbd.service
               └─4497 /home/dhpcsa/slurm-21.08.8/sbin/slurmdbd -D -s

Jan 28 12:57:40 controller slurmdbd[4497]: (null): _log_init: Unable to open logfile `/var/log/slurm/slurmdbd.log': No such file or directory
Jan 28 12:57:40 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
Jan 28 12:57:40 controller slurmdbd[4497]: slurmdbd: accounting_storage/as_mysql: _check_mysql_concat_is_sane: MySQL server version is: 5.5.5-10.3.39-M
Jan 28 12:57:40 controller slurmdbd[4497]: slurmdbd: error: Database settings not recommended values: innodb_buffer_pool_size innodb_lock_wait_timeout
Jan 28 12:57:40 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
Jan 28 12:57:41 controller slurmdbd[4497]: slurmdbd: slurmdbd version 21.08.8 started
Jan 28 12:57:41 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
Jan 28 12:57:41 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
Jan 28 12:57:46 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
Jan 28 13:00:01 controller slurmdbd[4497]: WARNING: MYSQL_OPT_RECONNECT is deprecated and will be removed in a future version.
[lines 1-19/19 (END)]
```

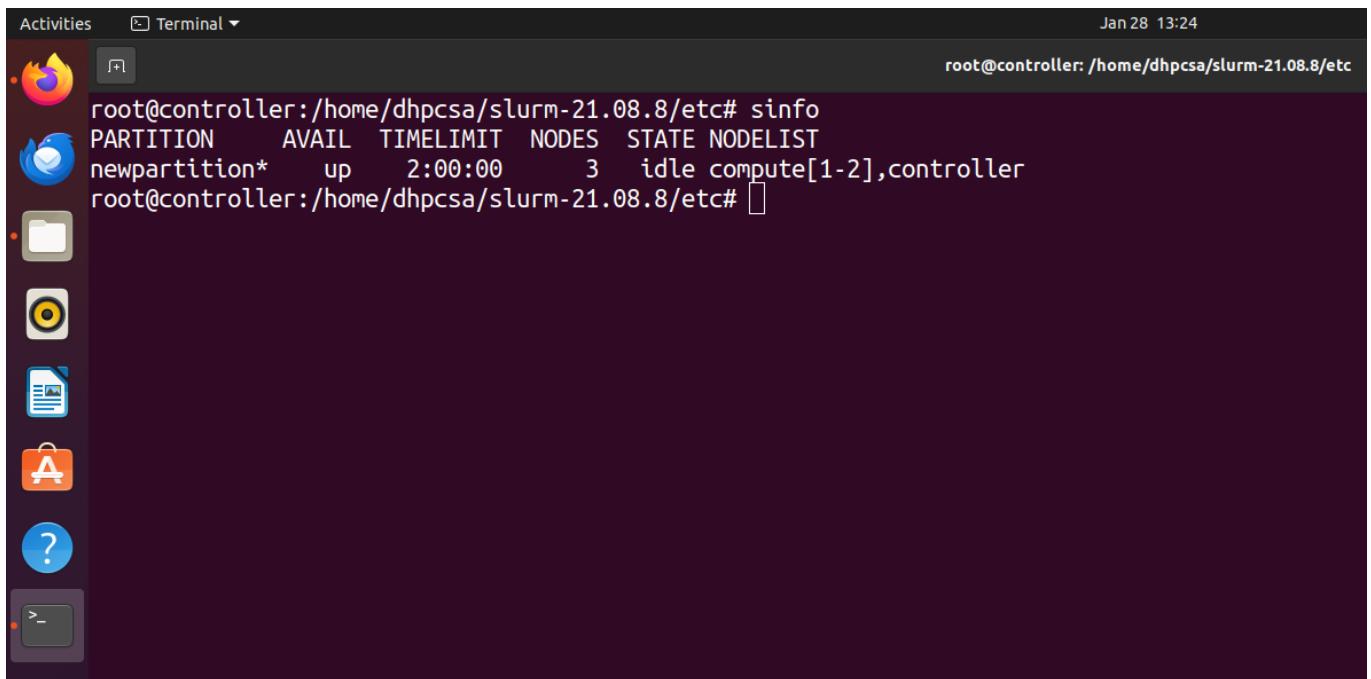
RUNNING STATUS OF MUNGE FOR AUTHENTICATION



```
root@controller:/home/dhpcsa/slurm-21.08.8/etc# systemctl status munge
● munge.service - MUNGE authentication service
  Loaded: loaded (/lib/systemd/system/munge.service; enabled; vendor preset: enabled)
  Active: active (running) since Tue 2025-01-28 12:30:04 IST; 52min ago
    Docs: man:munged(8)
   Main PID: 918 (munged)
     Tasks: 4 (limit: 5748)
    Memory: 1.6M
       CGroup: /system.slice/munge.service
               └─918 /usr/sbin/munged

Jan 28 12:30:03 controller systemd[1]: Starting MUNGE authentication service...
Jan 28 12:30:04 controller systemd[1]: Started MUNGE authentication service.
root@controller:/home/dhpcsa/slurm-21.08.8/etc#
```

DISPLAYS INFORMATION ABOUT THE STATE OF PARTITIONS, AND NODES



A screenshot of a Linux desktop environment, likely Ubuntu, showing a terminal window titled "Terminal". The terminal window displays the output of the "sinfo" command from the Slurm package. The output shows one partition named "newpartition*" with 3 nodes, all of which are currently "idle". The terminal window is located in the top right corner of the screen, and the desktop interface is visible on the left.

```
root@controller:/home/dhpcsa/slurm-21.08.8/etc# sinfo
PARTITION      AVAIL  TIMELIMIT  NODES  STATE NODELIST
newpartition*   up        2:00:00     3    idle  compute[1-2],controller
root@controller:/home/dhpcsa/slurm-21.08.8/etc#
```

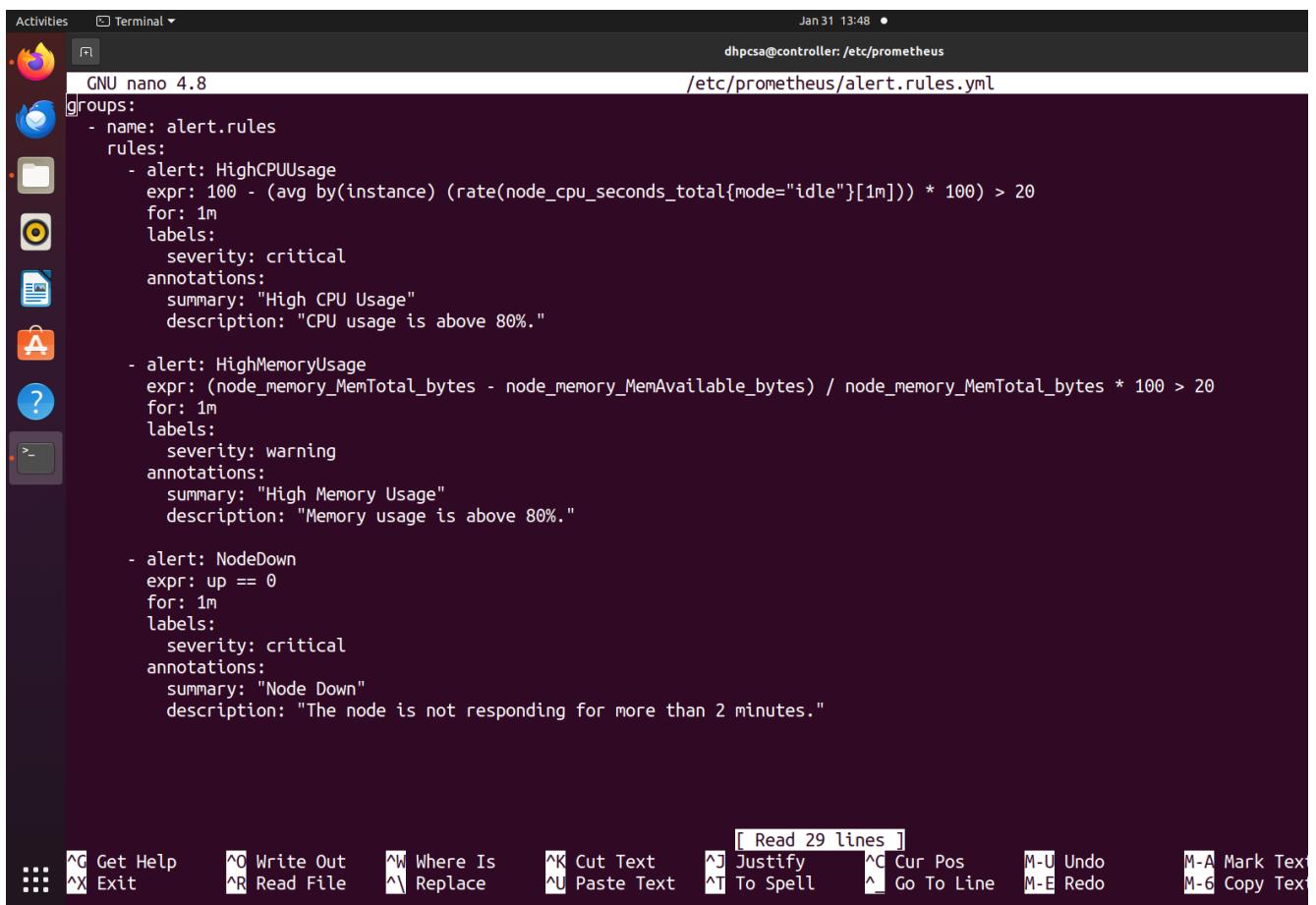
2. Monitoring Tools Installation

Description: This step involves setting up monitoring tools, including Prometheus, Grafana, Node Exporter, and Alertmanager. These tools will help collect and visualize metrics from the HPC cluster.

- **Prometheus:** A monitoring and alerting toolkit.
- **Grafana:** A visualization tool for monitoring data.
- **Node Exporter:** A Prometheus exporter for hardware and OS metrics.
- **Alertmanager:** Manages alerts sent by Prometheus.

Screenshot:

/etc/prometheus/alert.rules.yml File



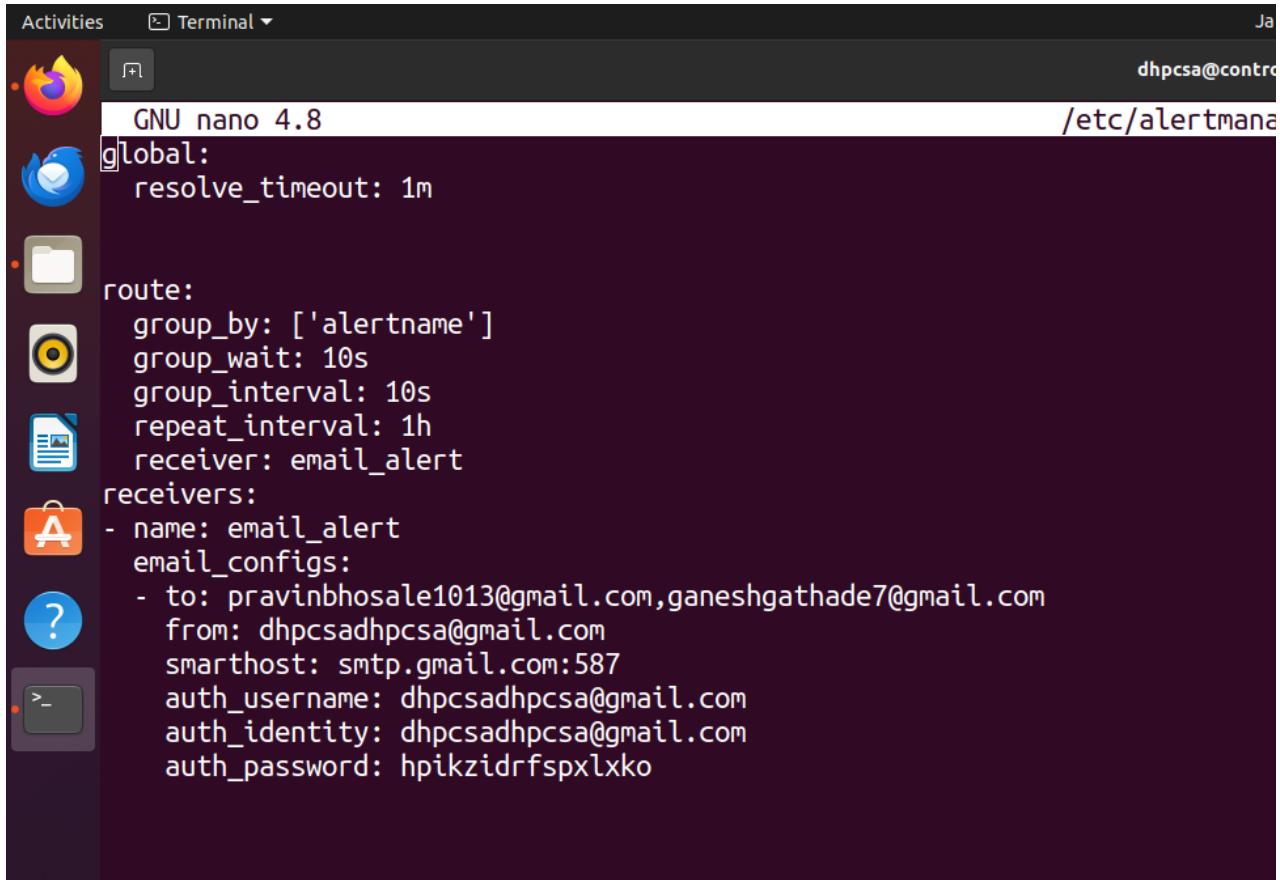
```

GNU nano 4.8
groups:
  - name: alert.rules
    rules:
      - alert: HighCPUUsage
        expr: 100 - (avg by(instance) (rate(node_cpu_seconds_total{mode="idle"}[1m])) * 100) > 20
        for: 1m
        labels:
          severity: critical
        annotations:
          summary: "High CPU Usage"
          description: "CPU usage is above 80%."
      - alert: HighMemoryUsage
        expr: (node_memory_MemTotal_bytes - node_memory_MemAvailable_bytes) / node_memory_MemTotal_bytes * 100 > 20
        for: 1m
        labels:
          severity: warning
        annotations:
          summary: "High Memory Usage"
          description: "Memory usage is above 80%."
      - alert: NodeDown
        expr: up == 0
        for: 1m
        labels:
          severity: critical
        annotations:
          summary: "Node Down"
          description: "The node is not responding for more than 2 minutes."

```

The terminal window shows the file content in a dark-themed nano editor. The file defines three alert groups: 'alert.rules'. Each group contains one or more alert definitions. Each alert has an 'expr' (expression), 'for' (duration), 'labels' (severity), and 'annotations' (summary and description). The 'HighCPUUsage' alert checks if CPU usage is above 80% over a minute. The 'HighMemoryUsage' alert checks if memory usage is above 80% over a minute. The 'NodeDown' alert checks if a node is not responding for more than 2 minutes.

/etc/alertmanager/alertmanager.yml File

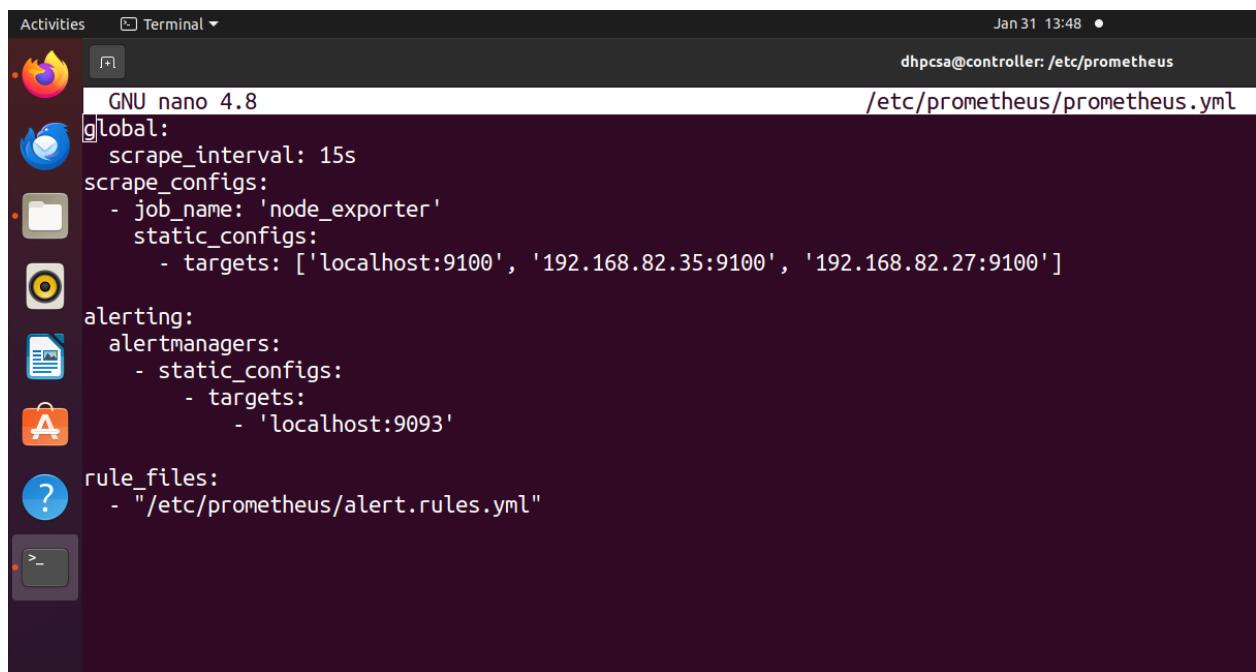


```
Activities Terminal ▾ Jan 31 13:48 dhpcsa@controller: /etc/alertmanager
GNU nano 4.8 /etc/alertmanager
global:
  resolve_timeout: 1m

route:
  group_by: ['alertname']
  group_wait: 10s
  group_interval: 10s
  repeat_interval: 1h
  receiver: email_alert

receivers:
- name: email_alert
  email_configs:
    - to: pravinbhosale1013@gmail.com,ganeshgathade7@gmail.com
      from: dhpcsa@dhpcsa@gmail.com
      smarthost: smtp.gmail.com:587
      auth_username: dhpcsa@dhpcsa@gmail.com
      auth_identity: dhpcsa@dhpcsa@gmail.com
      auth_password: hpikzidrfspxlxko
```

/etc/prometheus/prometheus.yml File

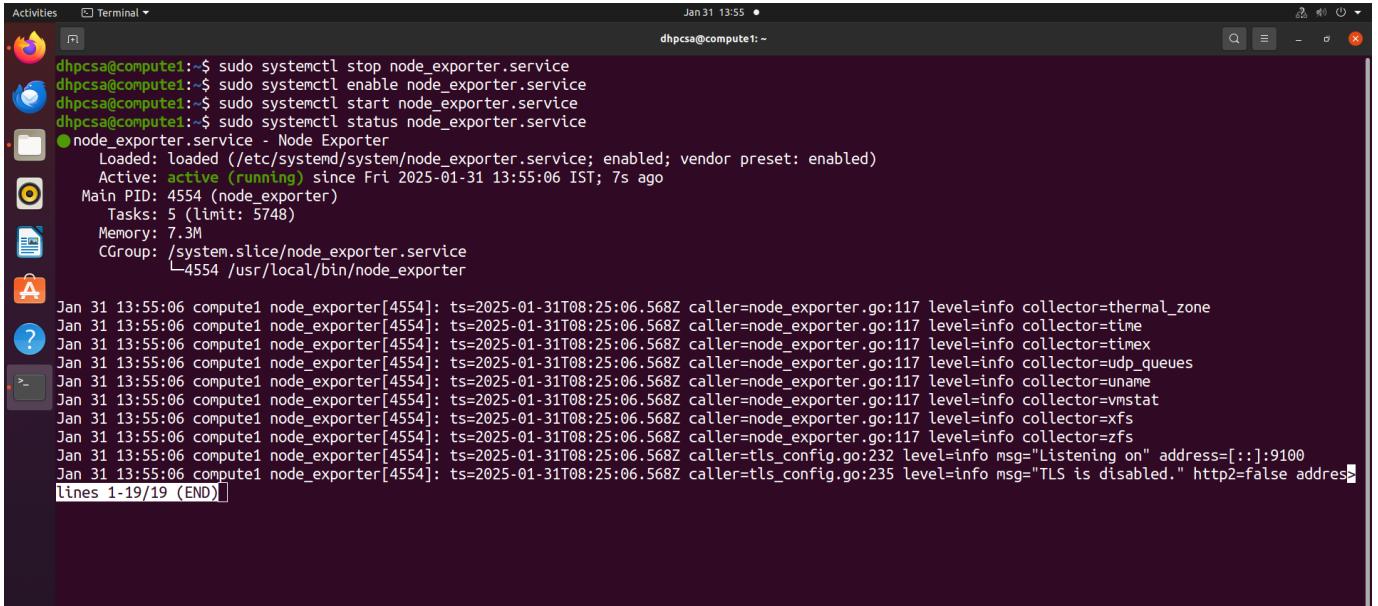


```
Activities Terminal ▾ Jan 31 13:48 dhpcsa@controller: /etc/prometheus
GNU nano 4.8 /etc/prometheus
global:
  scrape_interval: 15s
  scrape_configs:
    - job_name: 'node_exporter'
      static_configs:
        - targets: ['localhost:9100', '192.168.82.35:9100', '192.168.82.27:9100']

  alerting:
    alertmanagers:
      - static_configs:
          - targets:
              - 'localhost:9093'

  rule_files:
    - "/etc/prometheus/alert.rules.yml"
```

NODE EXPORTER STATUS



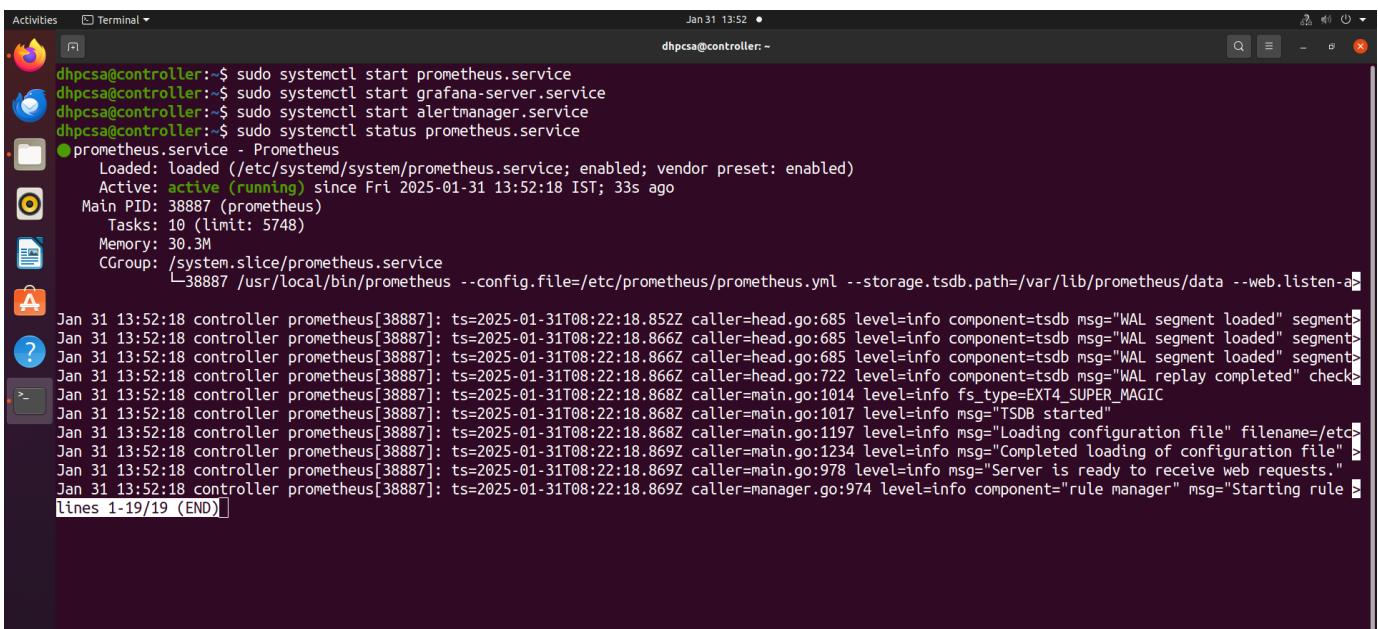
```

Activities Terminal Jan 31 13:55 • dhpcsa@compute1: ~
dhpcsa@compute1:~$ sudo systemctl stop node_exporter.service
dhpcsa@compute1:~$ sudo systemctl enable node_exporter.service
dhpcsa@compute1:~$ sudo systemctl start node_exporter.service
dhpcsa@compute1:~$ sudo systemctl status node_exporter.service
● node_exporter.service - Node Exporter
   Loaded: loaded (/etc/systemd/system/node_exporter.service; enabled; vendor preset: enabled)
   Active: active (running) since Fri 2025-01-31 13:55:06 IST; 7s ago
     Main PID: 4554 (node_exporter)
       Tasks: 5 (limit: 5748)
      Memory: 7.3M
        CGroup: /system.slice/node_exporter.service
                  └─4554 /usr/local/bin/node_exporter

Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=thermal_zone
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=time
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=timex
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=udp_queues
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=uname
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=vmstat
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=xfs
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=node_exporter.go:117 level=info collector=zfs
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=tls_config.go:232 level=info msg="Listening on" address=[::]:9100
Jan 31 13:55:06 compute1 node_exporter[4554]: ts=2025-01-31T08:25:06.568Z caller=tls_config.go:235 level=info msg="TLS is disabled." http2=false address=
lines 1-19/19 (END)

```

PROMETHEUS STATUS



```

Activities Terminal Jan 31 13:52 • dhpcsa@controller: ~
dhpcsa@controller:~$ sudo systemctl start prometheus.service
dhpcsa@controller:~$ sudo systemctl start grafana-server.service
dhpcsa@controller:~$ sudo systemctl start alertmanager.service
dhpcsa@controller:~$ sudo systemctl status prometheus.service
● prometheus.service - Prometheus
   Loaded: loaded (/etc/systemd/system/prometheus.service; enabled; vendor preset: enabled)
   Active: active (running) since Fri 2025-01-31 13:52:18 IST; 33s ago
     Main PID: 38887 (prometheus)
       Tasks: 10 (limit: 5748)
      Memory: 30.3M
        CGroup: /system.slice/prometheus.service
                  └─38887 /usr/local/bin/prometheus --config.file=/etc/prometheus/prometheus.yml --storage.tsdb.path=/var/lib/prometheus/data --web.listen-a

Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.852Z caller=head.go:685 level=info component=tsdb msg="WAL segment loaded" segment>
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.866Z caller=head.go:685 level=info component=tsdb msg="WAL segment loaded" segment>
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.866Z caller=head.go:685 level=info component=tsdb msg="WAL segment loaded" segment>
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.866Z caller=head.go:722 level=info component=tsdb msg="WAL replay completed" check>
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.868Z caller=main.go:1014 level=info fs_type=EXT4_SUPER_MAGIC
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.868Z caller=main.go:1017 level=info msg="TSDB started"
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.868Z caller=main.go:1197 level=info msg="Loading configuration file" filename=/etc>
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.869Z caller=main.go:1234 level=info msg="Completed loading of configuration file" >
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.869Z caller=main.go:978 level=info msg="Server is ready to receive web requests."
Jan 31 13:52:18 controller prometheus[38887]: ts=2025-01-31T08:22:18.869Z caller=manager.go:974 level=info component="rule manager" msg="Starting rule >
lines 1-19/19 (END)

```

GRAFANA STATUS

```
dhpcsa@controller:~$ sudo systemctl status grafana-server.service
● grafana-server.service - Grafana instance
   Loaded: loaded (/lib/systemd/system/grafana-server.service; enabled; vendor preset: enabled)
     Active: active (running) since Fri 2025-01-31 13:52:29 IST; 42s ago
       Docs: http://docs.grafana.org
      Main PID: 38938 (grafana)
        Tasks: 15 (limit: 5748)
       Memory: 53.6M
          CGroup: /system.slice/grafana-server.service
                  └─38938 /usr/share/grafana/bin/grafana server --config=/etc/grafana/grafana.ini --pidfile=/run/grafana/grafana-server.pid --packaging=deb

Jan 31 13:52:29 controller grafana[38938]: logger=ticker t=2025-01-31T13:52:29.601982656+05:30 level=info msg=starting first_tick=2025-01-31T13:52:30+05:30
Jan 31 13:52:29 controller grafana[38938]: logger=http.server t=2025-01-31T13:52:29.6032293+05:30 level=info msg="HTTP Server Listen" address=[::]:3000
Jan 31 13:52:29 controller grafana[38938]: logger=provisioning.dashboard t=2025-01-31T13:52:29.66729335+05:30 level=info msg="starting to provision das...
Jan 31 13:52:29 controller grafana[38938]: logger=provisioning.dashboard t=2025-01-31T13:52:29.667483134+05:30 level=info msg="finished to provision da...
Jan 31 13:52:29 controller grafana[38938]: logger=plugins.update.checker t=2025-01-31T13:52:29.926831804+05:30 level=info msg="Update check succeeded" >
Jan 31 13:52:29 controller grafana[38938]: logger=grafana.update.checker t=2025-01-31T13:52:29.958646295+05:30 level=info msg="Update check succeeded" >
Jan 31 13:52:30 controller grafana[38938]: logger=grafana.apiserver t=2025-01-31T13:52:30.119132638+05:30 level=info msg="Adding GroupVersion playlist...
Jan 31 13:52:30 controller grafana[38938]: logger=grafana.apiserver t=2025-01-31T13:52:30.119699162+05:30 level=info msg="Adding GroupVersion feature...
Jan 31 13:52:45 controller grafana[38938]: logger=context userId=1 orgId=1 uname=admin t=2025-01-31T13:52:45.381306276+05:30 level=info msg="Request Co...
Jan 31 13:52:45 controller grafana[38938]: logger=live t=2025-01-31T13:52:45.383769775+05:30 level=info msg="Initialized channel handler" channel=graf...
lines 1-20/20 (END)
```

ALERTMANAGER STATUS

```
dhpcsa@controller:~$ sudo systemctl status alertmanager.service
● alertmanager.service - Alertmanager
   Loaded: loaded (/etc/systemd/system/alertmanager.service; enabled; vendor preset: enabled)
     Active: active (running) since Fri 2025-01-31 13:52:35 IST; 50s ago
       Main PID: 38974 (alertmanager)
         Tasks: 10 (limit: 5748)
        Memory: 13.8M
          CGroup: /system.slice/alertmanager.service
                  └─38974 /usr/local/bin/alertmanager --config.file=/etc/alertmanager/alertmanager.yml --storage.path=/var/lib/alertmanager/data

Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.301Z caller=main.go:240 level=info msg="Starting Alertmanager" version="(version>
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.301Z caller=main.go:241 level=info build_context="(go=g01.19.4, user=root@abe866...
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.301Z caller=cluster.go:185 level=info component=cluster msg="Setting advertise a...
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.302Z caller=cluster.go:681 level=info component=cluster msg="Waiting for gossip >
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.328Z caller=coordinator.go:113 level=info component=configuration msg="Loading c...
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.328Z caller=coordinator.go:126 level=info component=configuration msg="Completed...
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.334Z caller=tls_config.go:232 level=info msg="Listening on" address=[::]:9093
Jan 31 13:52:35 controller alertmanager[38974]: ts=2025-01-31T08:22:35.334Z caller=tls_config.go:235 level=info msg="TLS is disabled." http2=false addr...
Jan 31 13:52:37 controller alertmanager[38974]: ts=2025-01-31T08:22:37.303Z caller=cluster.go:706 level=info component=cluster msg="gossip not settled">
Jan 31 13:52:45 controller alertmanager[38974]: ts=2025-01-31T08:22:45.307Z caller=cluster.go:698 level=info component=cluster msg="gossip settled; pro...
dhpcsa@controller:~$
```

DASHBOARD

ALERTMANAGER UI

The screenshot shows the Alertmanager UI interface. At the top, there's a navigation bar with tabs for Alertmanager, Alerts, Silences, Status, Settings, and Help. A "New Silence" button is located in the top right corner. Below the navigation, there's a search/filter section with tabs for "Filter" and "Group". A "Custom matcher, e.g. env='production'" input field is present. To the right, there are buttons for "Receiver: All", "Silenced", and "Inhibited". A "Silence" button is also visible. The main area displays a single alert card:

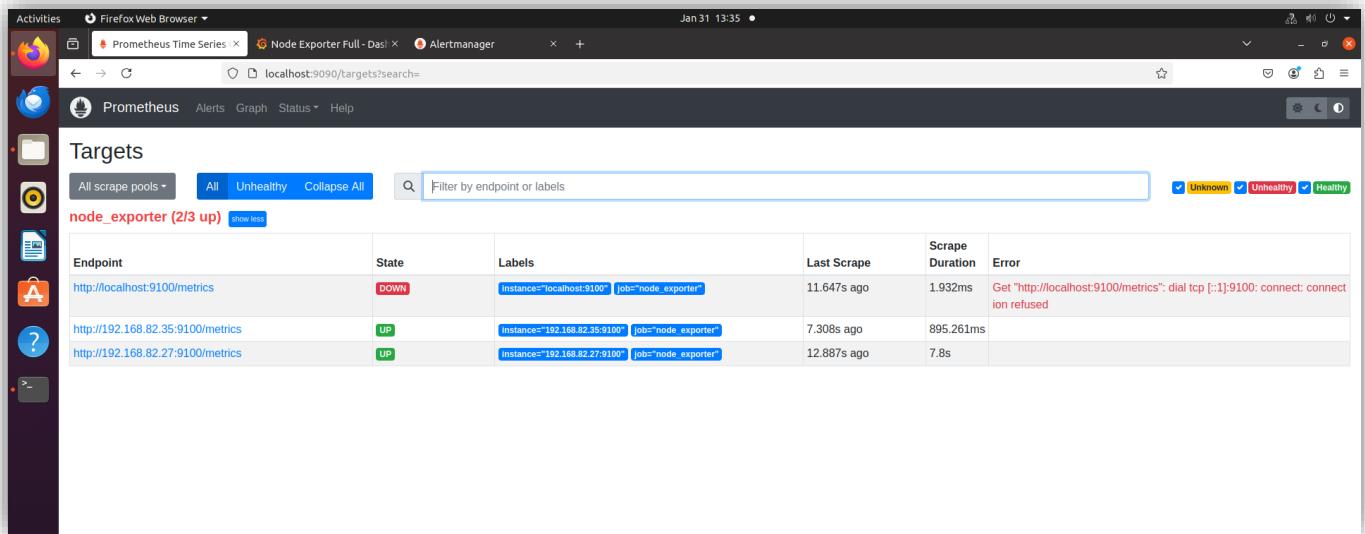
- Alertname:** NodeDown
- Timestamp:** 2025-01-31T07:55:46.014Z
- Labels:** instance="localhost:9100", job=node_exporter, severity=critical
- Annotations:** (empty)
- Actions:** Info, Source, Silence, Link

This screenshot shows the Alertmanager UI with two alerts listed:

- Alertname:** HighCPUUsage
- Alertname:** NodeDown

The interface is identical to the one in the first screenshot, with the same navigation, filter options, and alert details.

PROMETHEUS UI



The screenshot shows the Prometheus UI Targets page. The browser title bar indicates the URL is `localhost:9090/targets?search=node_exporter`. The page header includes links for Prometheus, Alerts, Graph, Status, and Help. A search bar at the top right allows filtering by endpoint or labels. Below the search bar, a button group shows the current filter: `All`, `Unhealthy`, and `Collapse All`. A checkbox group at the top right includes `Unknown`, `Unhealthy`, and `Healthy`. The main table displays three entries under the heading `node_exporter (2/3 up)`:

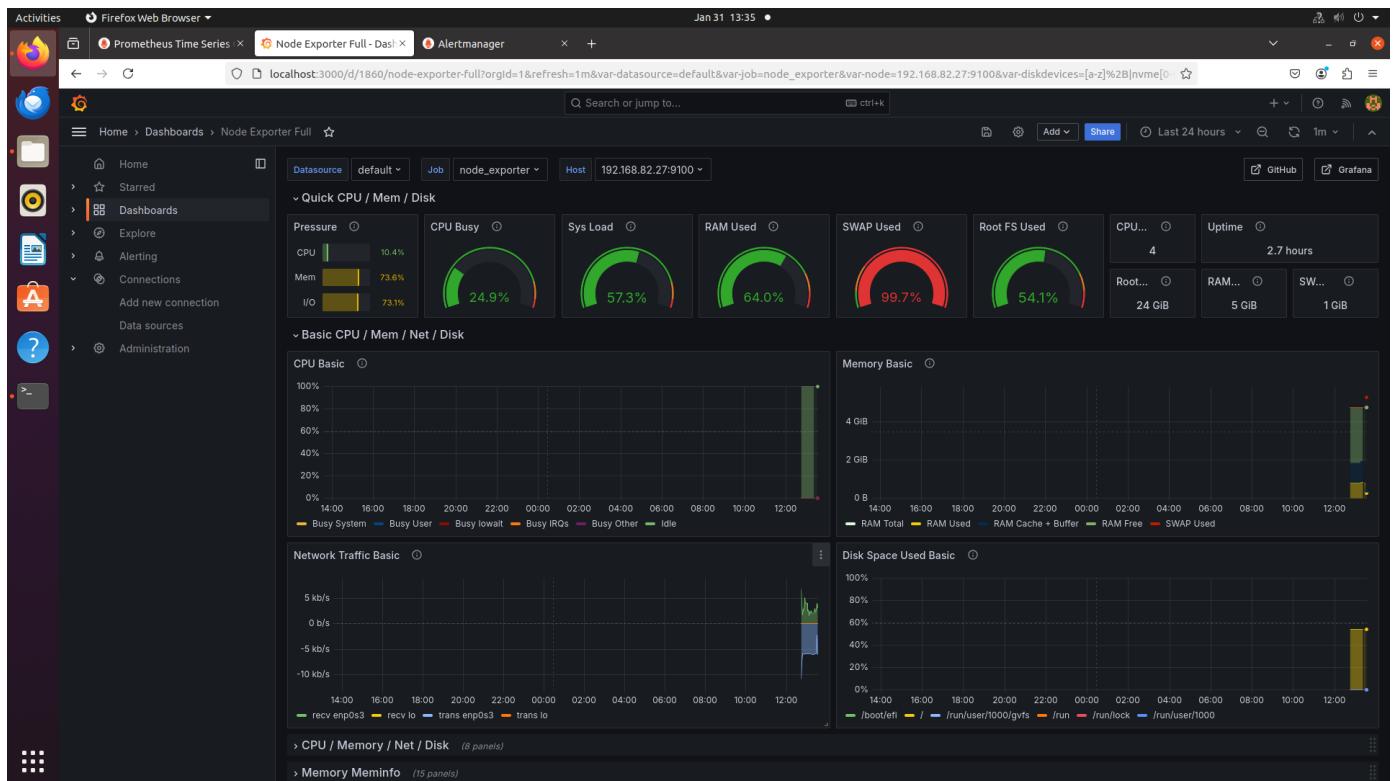
Endpoint	State	Labels	Last Scrape	Scrape Duration	Error
<code>http://localhost:9100/metrics</code>	DOWN	<code>instance="localhost:9100" job="node_exporter"</code>	11.647s ago	1.932ms	Get " <code>http://localhost:9100/metrics</code> ": dial tcp [::1]:9100: connect: connection refused
<code>http://192.168.82.35:9100/metrics</code>	UP	<code>instance="192.168.82.35:9100" job="node_exporter"</code>	7.308s ago	895.261ms	
<code>http://192.168.82.27:9100/metrics</code>	UP	<code>instance="192.168.82.27:9100" job="node_exporter"</code>	12.887s ago	7.8s	

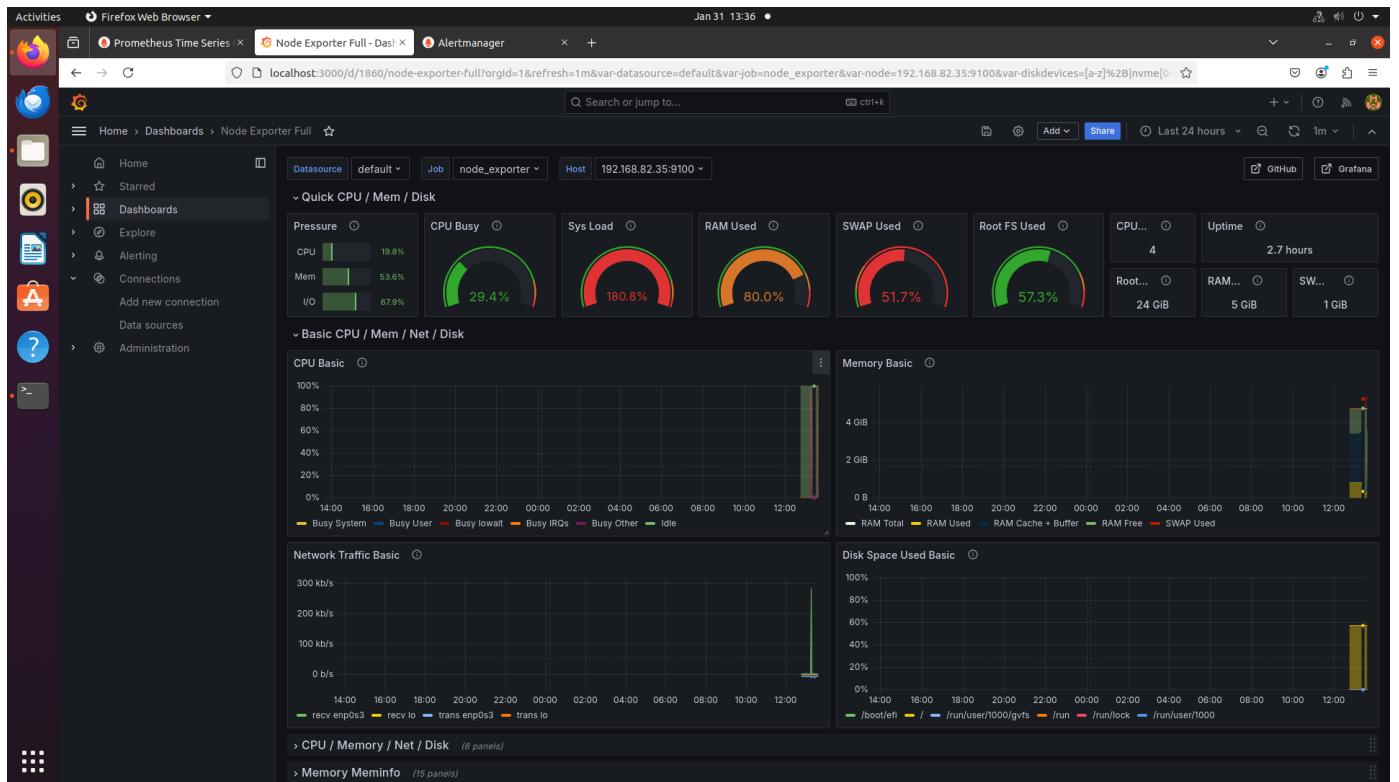
The screenshot shows the Prometheus Alertmanager interface with three active alerts listed:

- HighCPUUsage**: OK, 35.720s ago, 0.597ms evaluation time.
- HighMemoryUsage**: OK, 35.721s ago, 0.206ms evaluation time.
- NodeDown**: OK, 35.721s ago, 0.201ms evaluation time.

Each alert entry includes its configuration details, such as the alert name, expression, duration, labels, annotations, state, error, last evaluation time, and evaluation time.

GRAFANA DASHBOARD UI





EMAIL ALERTS

HighMemoryUsage Alert

2/8/25, 1:53 PM Gmail - [FIRING:3] HighMemoryUsage (node_exporter warning)

Gmail Pravin Bhosale <pravinbhosale1013@gmail.com>

[FIRING:3] HighMemoryUsage (node_exporter warning)

dhpcsdhpcsa@gmail.com <dhpcsdhpcsa@gmail.com>
To: pravinbhosale1013@gmail.com, ganeshgathade7@gmail.com 2 February 2025 at 14:35

3 alerts for alertname=HighMemoryUsage

[View In Alertmanager](#)

[3] Firing

Labels
alertname = HighMemoryUsage
instance = [192.168.82.27:9100](#)
job = node_exporter
severity = warning

Annotations
description = Memory usage is above 80%.
summary = High Memory Usage
[Source](#)

Labels
alertname = HighMemoryUsage
instance = [192.168.82.35:9100](#)
job = node_exporter
severity = warning

Annotations
description = Memory usage is above 80%.
summary = High Memory Usage
[Source](#)

NodeDown Alert

2/8/25, 2:02 PM Gmail - [FIRING:2] NodeDown (node_exporter critical)

Gmail Pravin Bhosale <pravinbhosale1013@gmail.com>

[FIRING:2] NodeDown (node_exporter critical)

dhpca@hpcsa@gmail.com <dhpca@hpcsa@gmail.com>
To: pravinbhosale1013@gmail.com, ganeshgathade7@gmail.com 31 January 2025 at 13:29

2 alerts for alertname=NodeDown

[View In Alertmanager](#)

[2] Firing

Labels
alertname = NodeDown
instance = 192.168.82.27:9100
job = node_exporter
severity = critical

Annotations
description = The node is not responding for more than 2 minutes.
summary = Node Down
[Source](#)

Labels
alertname = NodeDown
instance = localhost:9100
job = node_exporter
severity = critical

Annotations
description = The node is not responding for more than 2 minutes.
summary = Node Down
[Source](#)

Sent by Alertmanager

HighCPUUsage Alert

2/8/25, 2:09 PM

Gmail - [FIRING:2] HighCPUUsage (critical)



Pravin Bhosale <pravinbhosale1013@gmail.com>

[FIRING:2] HighCPUUsage (critical)

dhpcsdhpcsa@gmail.com <dhpcsdhpcsa@gmail.com>
To: pravinbhosale1013@gmail.com, ganeshgathade7@gmail.com

31 January 2025 at 13:36

2 alerts for alertname=HighCPUUsage

[View In Alertmanager](#)

[2] Firing

Labels

alertname = HighCPUUsage
instance = 192.168.82.27:9100
severity = critical

Annotations

description = CPU usage is above 80%.
summary = High CPU Usage

[Source](#)

Labels

alertname = HighCPUUsage
instance = 192.168.82.35:9100
severity = critical

Annotations

description = CPU usage is above 80%.
summary = High CPU Usage

[Source](#)

[Sent by Alertmanager](#)

3. Centralized Storage Configuration

Description: Set up NFS (Network File System) for centralized log storage. This allows all nodes in the HPC cluster to access and store logs in a shared location.

Screenshot

CENTRALIZED STORAGE: NFS (NETWORK FILE SYSTEM)

NFS SERVER SETUP SCRIPT

```

Activities Terminal ▾ Feb 9 12:39 •
grafana_11.2.1_amd64.deb      node_exporter-1.5.0.linux-amd64  prometheus-2.42.0.linux-amd64.tar.gz
dhpcsa@controller:~/project/nfs_setup$ cd nfs_setup/
dhpcsa@controller:~/project/nfs_setup$ ls
nfs_server.sh
dhpcsa@controller:~/project/nfs_setup$ cat nfs_server.sh
#!/bin/bash

# Variables
EXPORT_DIR="/mnt/nfs_project"
SLAVE_NODES=("192.168.82.35" "192.168.82.27") # Replace with actual slave IPs

# Install required packages
echo "Installing NFS server..."
sudo apt update -y
sudo apt install -y nfs-kernel-server

# Create and configure the shared directory
echo "Setting up shared directory..."
sudo mkdir -p $EXPORT_DIR
sudo chown nobody:nogroup $EXPORT_DIR
sudo chmod 755 $EXPORT_DIR

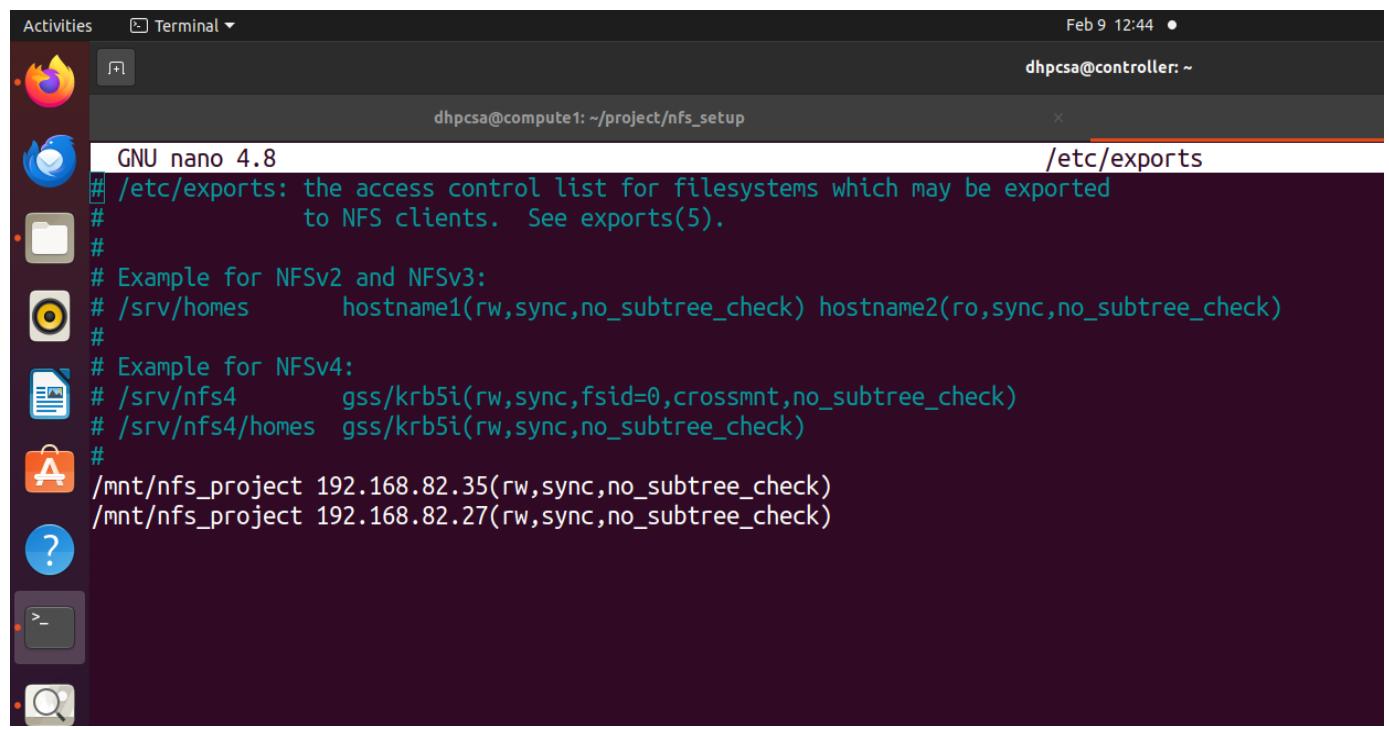
# Add export rules for the specified slave nodes
echo "Configuring export rules..."
for IP in "${SLAVE_NODES[@]}"; do
    echo "$EXPORT_DIR $IP(rw,sync,no_subtree_check)" | sudo tee -a /etc/exports
done

# Restart and enable NFS server
echo "Restarting NFS server..."
sudo systemctl restart nfs-kernel-server
sudo systemctl enable nfs-kernel-server

# Verify exports
sudo exportfs -v

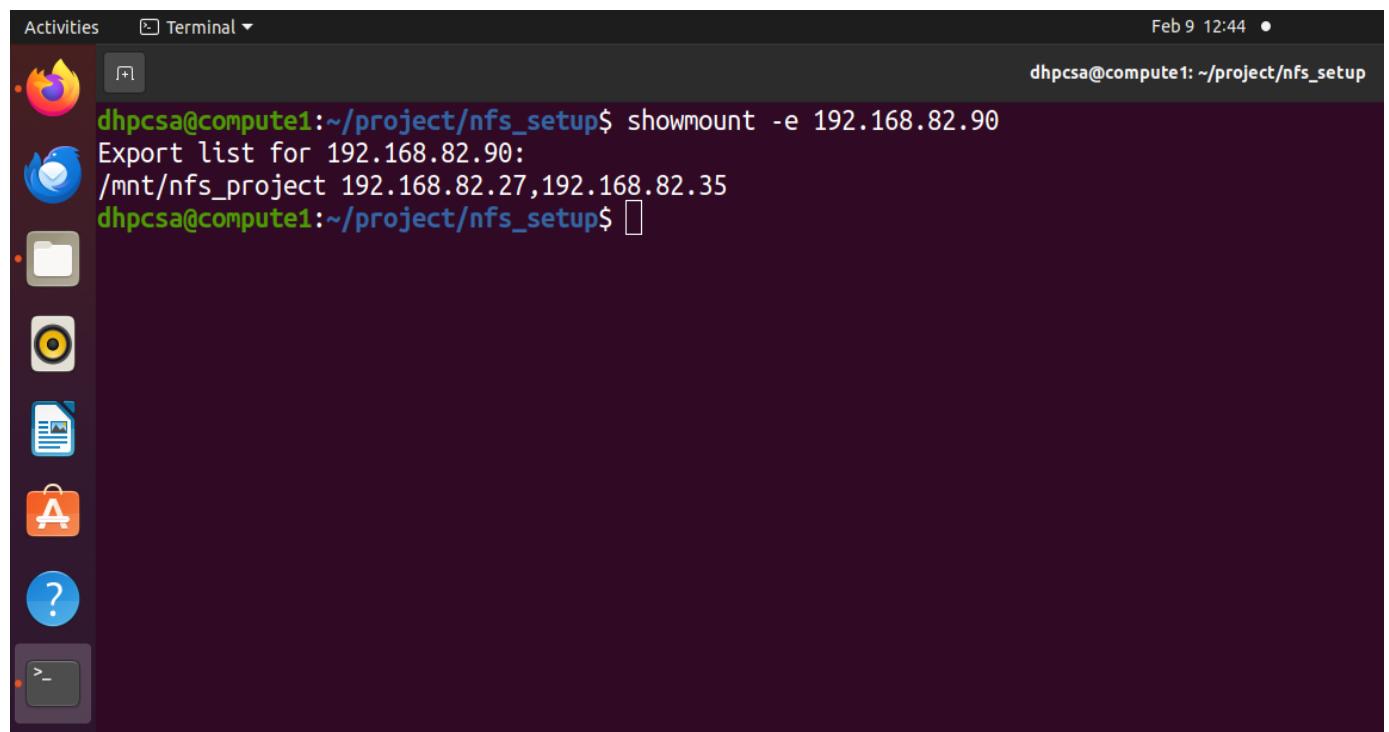
echo "NFS server setup completed!"
dhpcsa@controller:~/project/nfs_setup$ 

```



A screenshot of a Linux desktop environment, likely Ubuntu. On the left is a vertical dock with icons for various applications: a browser, a file manager, a terminal, a file browser, a system settings icon, a help icon, a terminal icon, and a search icon. The main window is a terminal titled "GNU nano 4.8" with the file "/etc/exports" open. The content of the file is:

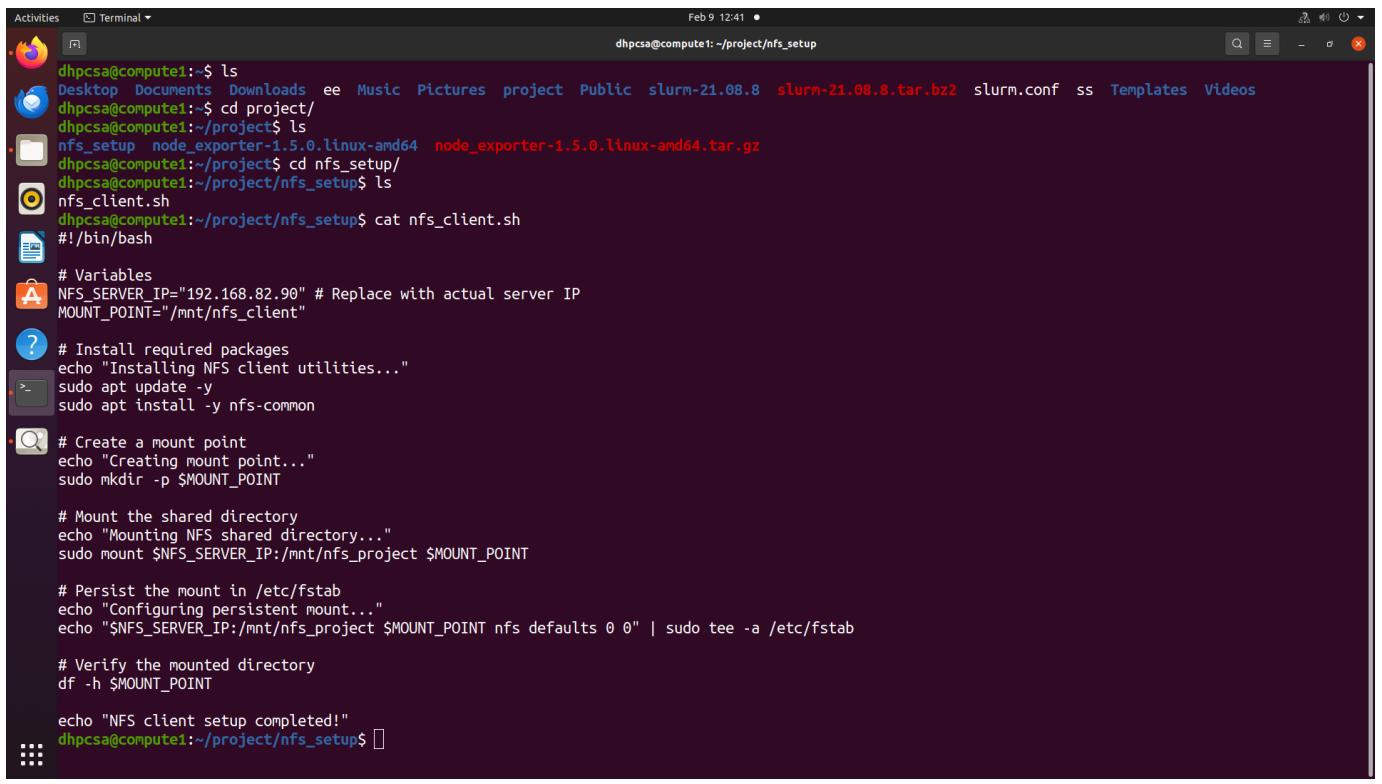
```
# /etc/exports: the access control list for filesystems which may be exported
#           to NFS clients. See exports(5).
#
# Example for NFSv2 and NFSv3:
# /srv/homes    hostname1(rw,sync,no_subtree_check) hostname2(ro,sync,no_subtree_check)
#
# Example for NFSv4:
# /srv/nfs4      gss/krb5i(rw,sync,fsid=0,crossmnt,no_subtree_check)
# /srv/nfs4/homes gss/krb5i(rw,sync,no_subtree_check)
#
# /mnt/nfs_project 192.168.82.35(rw,sync,no_subtree_check)
# /mnt/nfs_project 192.168.82.27(rw,sync,no_subtree_check)
```



A screenshot of a Linux desktop environment, likely Ubuntu. On the left is a vertical dock with icons for various applications: a browser, a file manager, a terminal, a file browser, a system settings icon, a help icon, a terminal icon, and a search icon. The main window is a terminal window with the command "showmount -e 192.168.82.90" run. The output is:

```
dhpcsa@compute1:~/project/nfs_setup$ showmount -e 192.168.82.90
Export list for 192.168.82.90:
/mnt/nfs_project 192.168.82.27,192.168.82.35
dhpcsa@compute1:~/project/nfs_setup$
```

NSF CLIENT SETUP SCRIPT



The screenshot shows a terminal window titled "Terminal" with the command "dhpcsa@compute1:~\$". The window displays a series of shell commands being run to set up an NFS client. The commands include navigating to a project directory, extracting files, and executing a script named "nfs_client.sh". The terminal also shows variable assignments for the NFS server IP and mount point, package installations, mount creation, directory mounting, persistent mount configuration, and a final verification step.

```
dhpcsa@compute1:~$ ls
Desktop Documents Downloads ee Music Pictures project Public slurm-21.08.8 slurm-21.08.8.tar.bz2 slurm.conf ss Templates Videos
dhpcsa@compute1:~$ cd project/
dhpcsa@compute1:~/project$ ls
nfs_setup node_exporter-1.5.0.linux-amd64 node_exporter-1.5.0.linux-amd64.tar.gz
dhpcsa@compute1:~/project$ cd nfs_setup/
dhpcsa@compute1:~/project/nfs_setup$ ls
nfs_client.sh
dhpcsa@compute1:~/project/nfs_setup$ cat nfs_client.sh
#!/bin/bash

# Variables
NFS_SERVER_IP="192.168.82.90" # Replace with actual server IP
MOUNT_POINT="/mnt/nfs_client"

# Install required packages
echo "Installing NFS client utilities..."
sudo apt update -y
sudo apt install -y nfs-common

# Create a mount point
echo "Creating mount point..."
sudo mkdir -p $MOUNT_POINT

# Mount the shared directory
echo "Mounting NFS shared directory..."
sudo mount $NFS_SERVER_IP:$mnt/nfs_project $MOUNT_POINT

# Persist the mount in /etc/fstab
echo "Configuring persistent mount..."
echo "$NFS_SERVER_IP:$mnt/nfs_project $MOUNT_POINT nfs defaults 0 0" | sudo tee -a /etc/fstab

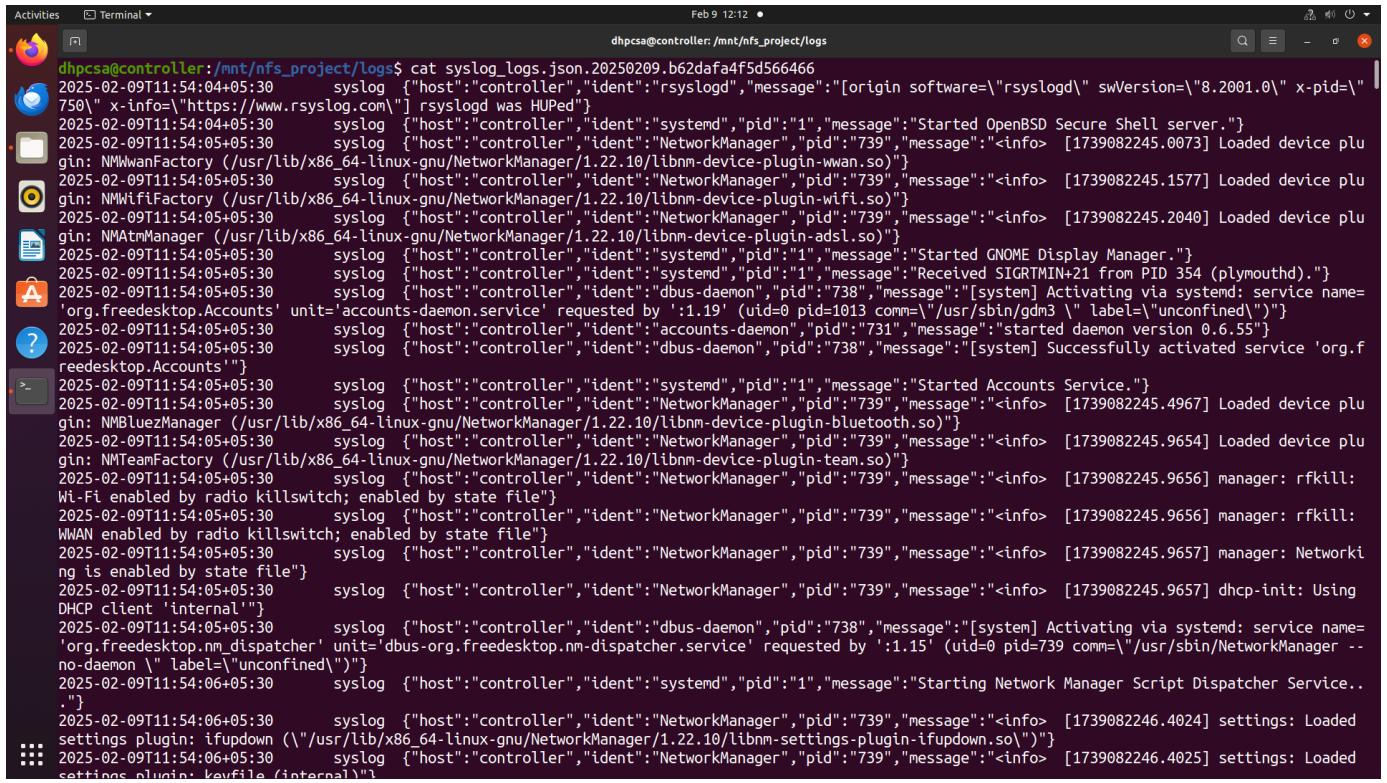
# Verify the mounted directory
df -h $MOUNT_POINT

echo "NFS client setup completed!"
dhpcsa@compute1:~/project/nfs_setup$
```

4. Log Parsing and Exploratory Data Analysis (EDA) for ML Model Training

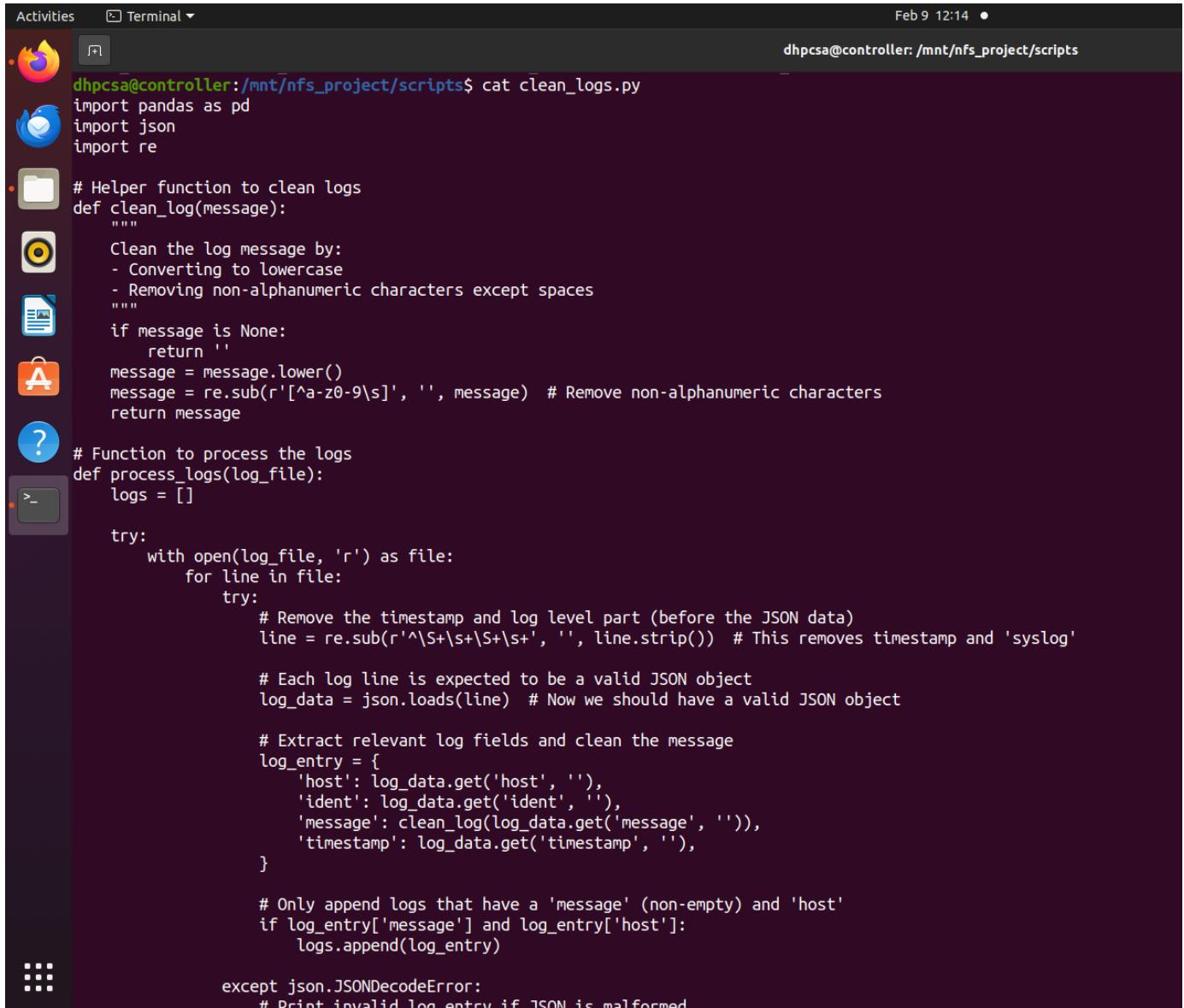
Description: Use Fluentd for log parsing and preprocessing. This step involves extracting relevant information from logs and preparing the data for machine learning model training.

CENTRALIZED LOG COLLECTION



```
dhpcsa@controller:/mnt/nfs_project/logs$ cat syslog_logs.json.20250209.b62dafa4f5d566466
2025-02-09T11:54:04+05:30    syslog {"host":"controller","ident":"rsyslogd","message":"[origin software=\"rsyslogd\" swVersion=\"8.2001.0\" x-pid=\"750\" x-info=\"https://www.rsyslog.com/] rsyslogd was HUPed"}
2025-02-09T11:54:04+05:30    syslog {"host":"controller","ident":"systemd","pid":1,"message":"Started OpenBSD Secure Shell server."}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.0073] Loaded device plugin: NMWwanFactory (/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-device-plugin-wwan.so)"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.1577] Loaded device plugin: NMWifiFactory (/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-device-plugin-wifi.so)"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.2040] Loaded device plugin: NMAtmManager (/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-device-plugin-adsl.so)"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"systemd","pid":1,"message":"Started GNOME Display Manager."}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"systemd","pid":1,"message":"Received SIGRTMIN+21 from PID 354 (plymouthd)."}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"dbus-daemon","pid":738,"message":">[system] Activating via systemd: service name='org.freedesktop.Accounts' unit='accounts-daemon.service' requested by '1:19' (uid=0 pid=1013 comm="/usr/sbin/gdm3 \\" label=\"unconfined\")"
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"accounts-daemon","pid":731,"message":">started daemon version 0.6.55"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"dbus-daemon","pid":738,"message">[system] Successfully activated service 'org.freedesktop.Accounts'"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"systemd","pid":1,"message":"Started Accounts Service."}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.4967] Loaded device plugin: NMBluezManager (/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-device-plugin-bluetooth.so)"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.9654] Loaded device plugin: NMTeamFactory (/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-device-plugin-team.so)"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.9656] manager: rfkill: Wi-Fi enabled by radio killswitch; enabled by state file"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.9656] manager: rfkill: WWAN enabled by radio killswitch; enabled by state file"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.9657] manager: Networkking is enabled by state file"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082245.9657] dhcp-init: Using DHCP client 'internal'"}
2025-02-09T11:54:05+05:30    syslog {"host":"controller","ident":"dbus-daemon","pid":738,"message">[system] Activating via systemd: service name='org.freedesktop.nm_dispatcher' unit='dbus-org.freedesktop.nm-dispatcher.service' requested by '1:15' (uid=0 pid=739 comm="/usr/sbin/NetworkManager --no-daemon \\" label=\"unconfined\")"
2025-02-09T11:54:06+05:30    syslog {"host":"controller","ident":"systemd","pid":1,"message">Starting Network Manager Script Dispatcher Service.. ."}
2025-02-09T11:54:06+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082246.4024] settings: Loaded settings plugin: ifupdown ("(/usr/lib/x86_64-linux-gnu/NetworkManager/1.22.10/libnm-settings-plugin-ifupdown.so)")"
2025-02-09T11:54:06+05:30    syslog {"host":"controller","ident":"NetworkManager","pid":739,"message":"><info> [1739082246.4025] settings: Loaded settings plugin: keyfile (internal)"}
```

SCRIPT TO CLEAN LOGS



```

Activities Terminal ▾ Feb 9 12:14 •
dhpcsa@controller:/mnt/nfs_project/scripts$ cat clean_logs.py
import pandas as pd
import json
import re

# Helper function to clean logs
def clean_log(message):
    """
    Clean the log message by:
    - Converting to lowercase
    - Removing non-alphanumeric characters except spaces
    """
    if message is None:
        return ''
    message = message.lower()
    message = re.sub(r'^[^\w\s]+|\s+|[\w\s]+$', '', message) # Remove non-alphanumeric characters
    return message

# Function to process the logs
def process_logs(log_file):
    logs = []

    try:
        with open(log_file, 'r') as file:
            for line in file:
                try:
                    # Remove the timestamp and log level part (before the JSON data)
                    line = re.sub(r'^\S+\s+\S+\s+', '', line.strip()) # This removes timestamp and 'syslog'

                    # Each log line is expected to be a valid JSON object
                    log_data = json.loads(line) # Now we should have a valid JSON object

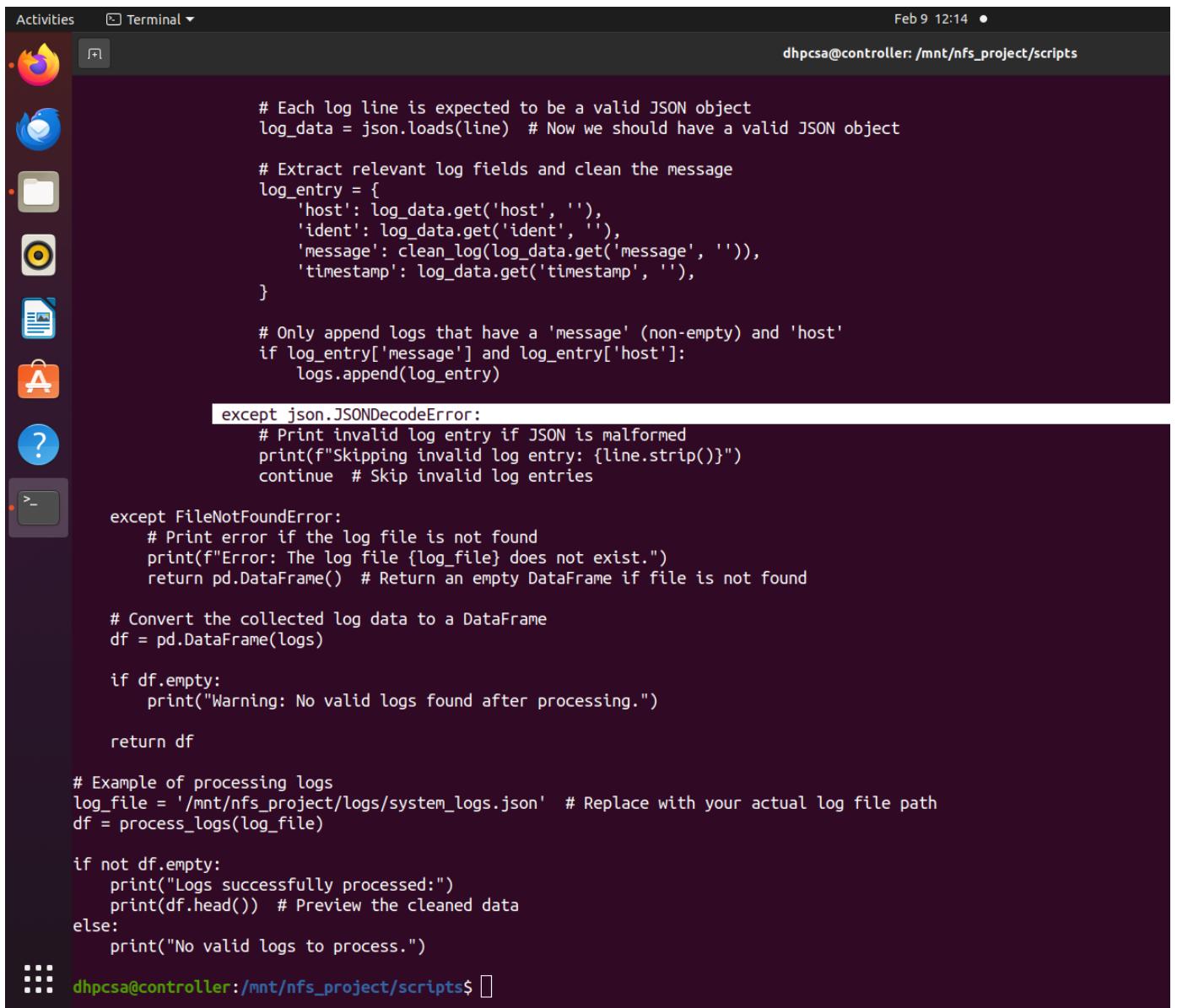
                    # Extract relevant log fields and clean the message
                    log_entry = {
                        'host': log_data.get('host', ''),
                        'ident': log_data.get('ident', ''),
                        'message': clean_log(log_data.get('message', '')),
                        'timestamp': log_data.get('timestamp', ''),
                    }

                    # Only append logs that have a 'message' (non-empty) and 'host'
                    if log_entry['message'] and log_entry['host']:
                        logs.append(log_entry)

                except json.JSONDecodeError:
                    # Print invalid log entry if JSON is malformed
                    print(f"Invalid JSON entry: {line}")
    except FileNotFoundError:
        print(f"File '{log_file}' not found")
    except Exception as e:
        print(f"An error occurred: {e}")

    return logs

```



The screenshot shows a Linux desktop environment with a terminal window open. The terminal window has a dark background and contains Python code for processing logs. The code uses json.loads to parse log lines into JSON objects, then extracts relevant fields like host, ident, message, and timestamp. It filters out logs without a message or host. It handles JSONDecodeError by skipping invalid entries. If a log file is not found, it prints an error and returns an empty DataFrame. Finally, it converts the collected logs into a DataFrame and prints a success message if it's not empty, or a warning if no logs were found.

```
# Each log line is expected to be a valid JSON object
log_data = json.loads(line) # Now we should have a valid JSON object

# Extract relevant log fields and clean the message
log_entry = {
    'host': log_data.get('host', ''),
    'ident': log_data.get('ident', ''),
    'message': clean_log(log_data.get('message', '')),
    'timestamp': log_data.get('timestamp', ''),
}

# Only append logs that have a 'message' (non-empty) and 'host'
if log_entry['message'] and log_entry['host']:
    logs.append(log_entry)

except json.JSONDecodeError:
    # Print invalid log entry if JSON is malformed
    print(f"Skipping invalid log entry: {line.strip()}")
    continue # Skip invalid log entries

except FileNotFoundError:
    # Print error if the log file is not found
    print(f"Error: The log file {log_file} does not exist.")
    return pd.DataFrame() # Return an empty DataFrame if file is not found

# Convert the collected log data to a DataFrame
df = pd.DataFrame(logs)

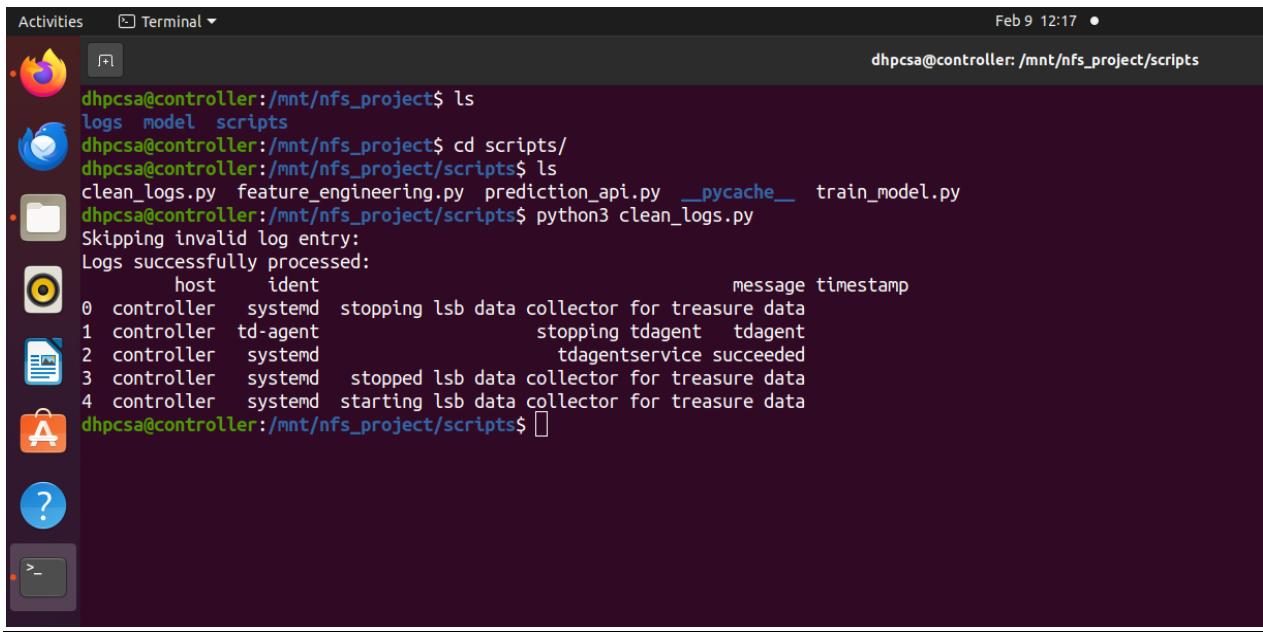
if df.empty:
    print("Warning: No valid logs found after processing.")

return df

# Example of processing logs
log_file = '/mnt/nfs_project/logs/system_logs.json' # Replace with your actual log file path
df = process_logs(log_file)

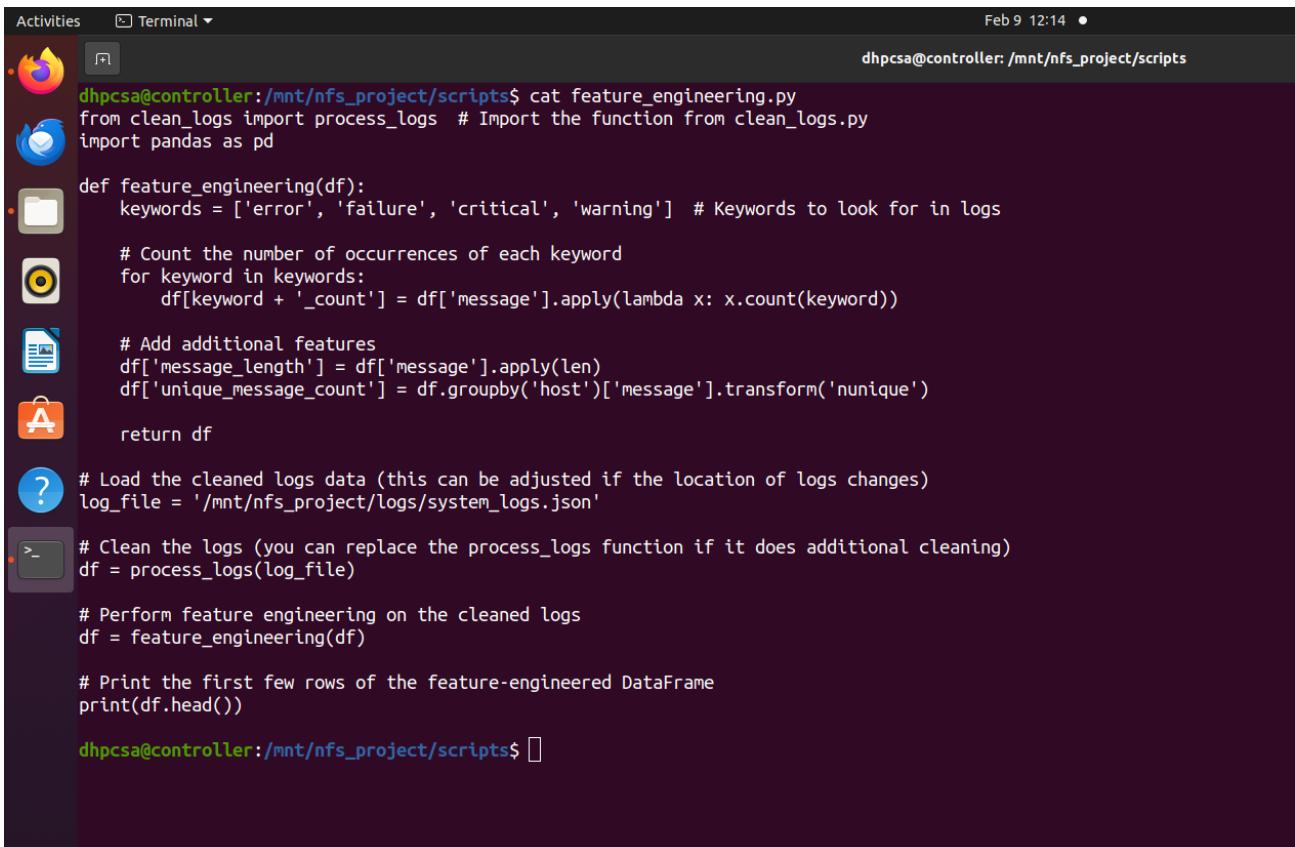
if not df.empty:
    print("Logs successfully processed:")
    print(df.head()) # Preview the cleaned data
else:
    print("No valid logs to process.")
```

RUN SCRIPT



```
dhpcsa@controller:/mnt/nfs_project$ ls
logs model scripts
dhpcsa@controller:/mnt/nfs_project$ cd scripts/
dhpcsa@controller:/mnt/nfs_project/scripts$ ls
clean_logs.py feature_engineering.py prediction_api.py __pycache__ train_model.py
dhpcsa@controller:/mnt/nfs_project/scripts$ python3 clean_logs.py
Skipping invalid log entry:
Logs successfully processed:
host ident message timestamp
0 controller systemd stopping lsb data collector for treasure data
1 controller td-agent stopping tdagent tdagent
2 controller systemd tdagentservice succeeded
3 controller systemd stopped lsb data collector for treasure data
4 controller systemd starting lsb data collector for treasure data
dhpcsa@controller:/mnt/nfs_project/scripts$ 
```

SCRIPT FOR FEATURE ENGINEERING



```
dhpcsa@controller:/mnt/nfs_project/scripts$ cat feature_engineering.py
from clean_logs import process_logs # Import the function from clean_logs.py
import pandas as pd

def feature_engineering(df):
    keywords = ['error', 'failure', 'critical', 'warning'] # Keywords to look for in logs

    # Count the number of occurrences of each keyword
    for keyword in keywords:
        df[keyword + '_count'] = df['message'].apply(lambda x: x.count(keyword))

    # Add additional features
    df['message_length'] = df['message'].apply(len)
    df['unique_message_count'] = df.groupby('host')['message'].transform('nunique')

    return df

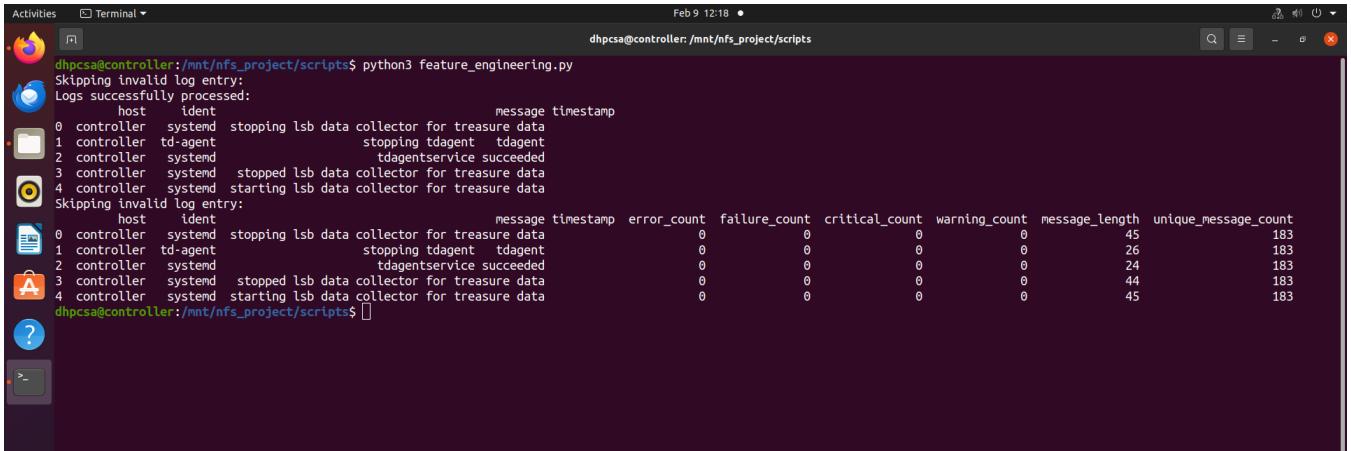
# Load the cleaned logs data (this can be adjusted if the location of logs changes)
log_file = '/mnt/nfs_project/logs/system_logs.json'

# Clean the logs (you can replace the process_logs function if it does additional cleaning)
df = process_logs(log_file)

# Perform feature engineering on the cleaned logs
df = feature_engineering(df)

# Print the first few rows of the feature-engineered DataFrame
print(df.head())
dhpcsa@controller:/mnt/nfs_project/scripts$ 
```

RUN SCRIPT



```

dhpcsa@controller:/mnt/nfs_project/scripts$ python3 feature_engineering.py
Skipping invalid log entry:
Logs successfully processed:
host      ident          message timestamp
0 controller systemd stopping lsb data collector for treasure data
1 controller td-agent   stopping tdagent  tdagent
2 controller systemd   tdagentservice succeeded
3 controller systemd stopped lsb data collector for treasure data
4 controller systemd starting lsb data collector for treasure data
Skipping invalid log entry:
host      ident          message timestamp error_count failure_count critical_count warning_count message_length unique_message_count
0 controller systemd stopping lsb data collector for treasure data 0 0 0 0 45 183
1 controller td-agent   stopping tdagent  tdagent 0 0 0 0 26 183
2 controller systemd   tdagentservice succeeded 0 0 0 0 24 183
3 controller systemd stopped lsb data collector for treasure data 0 0 0 0 44 183
4 controller systemd starting lsb data collector for treasure data 0 0 0 0 45 183
dhpcsa@controller:/mnt/nfs_project/scripts$ 
```

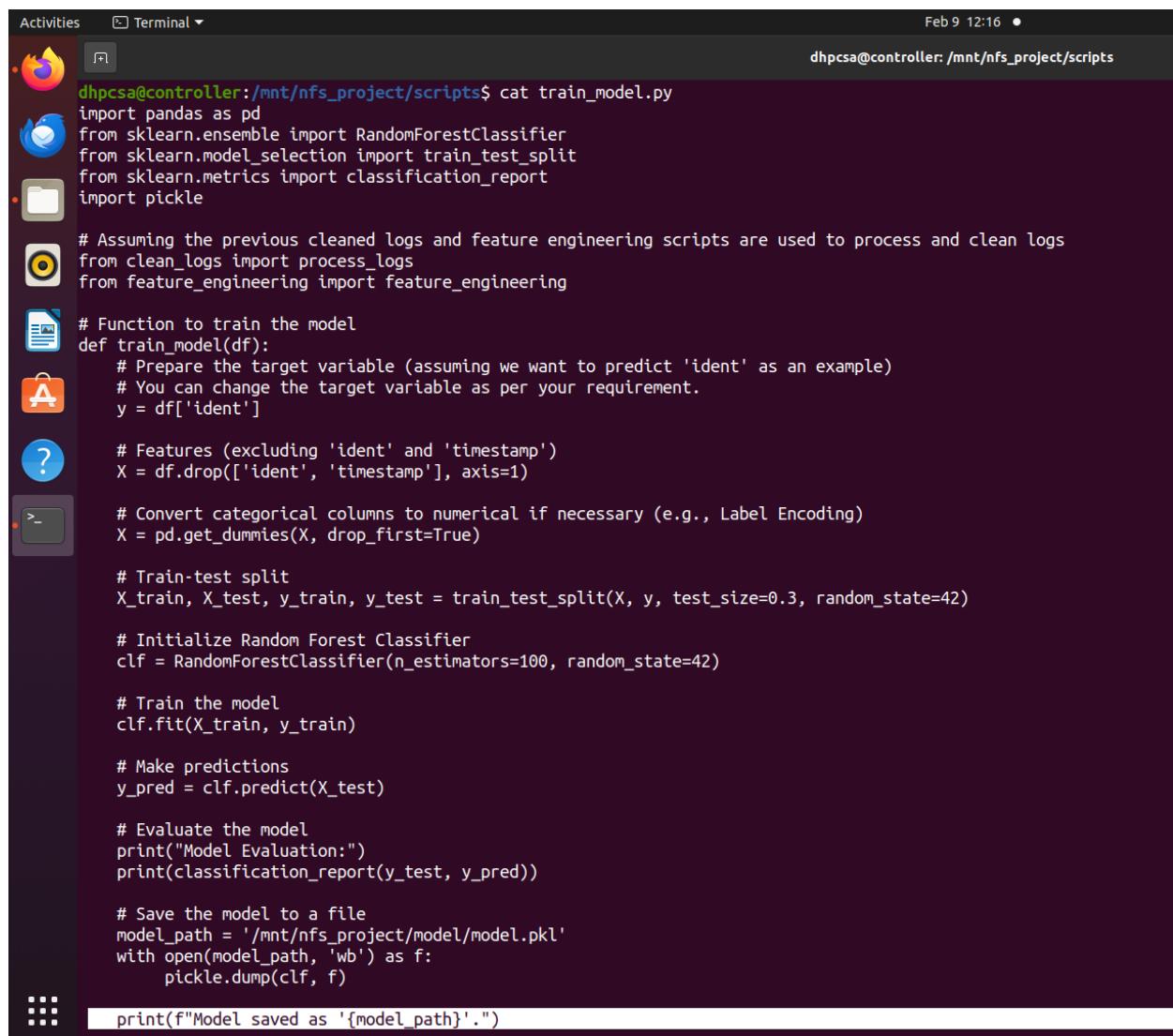
MACHINE LEARNING MODEL DEVELOPMENT

5. Model Development

Description: In this step, various machine learning models are developed and trained using the preprocessed data. This includes anomaly detection and predictive modeling techniques.

Screenshot:

MODEL TRAINING



```
dhpcsa@controller:/mnt/nfs_project/scripts$ cat train_model.py
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
import pickle

# Assuming the previous cleaned logs and feature engineering scripts are used to process and clean logs
from clean_logs import process_logs
from feature_engineering import feature_engineering

# Function to train the model
def train_model(df):
    # Prepare the target variable (assuming we want to predict 'ident' as an example)
    # You can change the target variable as per your requirement.
    y = df['ident']

    # Features (excluding 'ident' and 'timestamp')
    X = df.drop(['ident', 'timestamp'], axis=1)

    # Convert categorical columns to numerical if necessary (e.g., Label Encoding)
    X = pd.get_dummies(X, drop_first=True)

    # Train-test split
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

    # Initialize Random Forest Classifier
    clf = RandomForestClassifier(n_estimators=100, random_state=42)

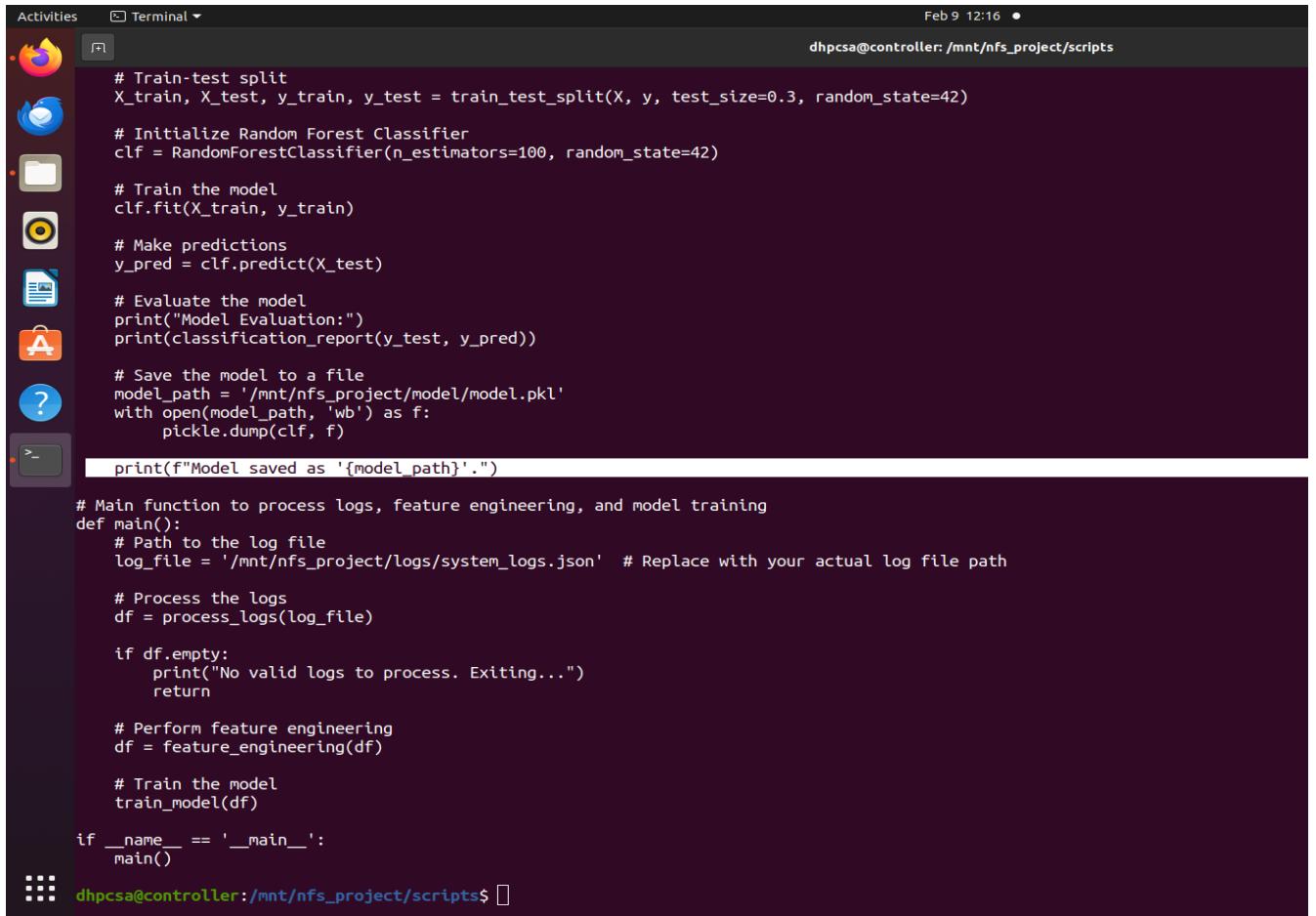
    # Train the model
    clf.fit(X_train, y_train)

    # Make predictions
    y_pred = clf.predict(X_test)

    # Evaluate the model
    print("Model Evaluation:")
    print(classification_report(y_test, y_pred))

    # Save the model to a file
    model_path = '/mnt/nfs_project/model/model.pkl'
    with open(model_path, 'wb') as f:
        pickle.dump(clf, f)

    print(f"Model saved as '{model_path}' .")
```



```

Activities Terminal Feb 9 12:16 •
dhpcsa@controller:/mnt/nfs_project/scripts

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize Random Forest Classifier
clf = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the model
clf.fit(X_train, y_train)

# Make predictions
y_pred = clf.predict(X_test)

# Evaluate the model
print("Model Evaluation:")
print(classification_report(y_test, y_pred))

# Save the model to a file
model_path = '/mnt/nfs_project/model/model.pkl'
with open(model_path, 'wb') as f:
    pickle.dump(clf, f)

print(f"Model saved as '{model_path}'.")

# Main function to process logs, feature engineering, and model training
def main():
    # Path to the log file
    log_file = '/mnt/nfs_project/logs/system_logs.json' # Replace with your actual log file path

    # Process the logs
    df = process_logs(log_file)

    if df.empty:
        print("No valid logs to process. Exiting...")
        return

    # Perform feature engineering
    df = feature_engineering(df)

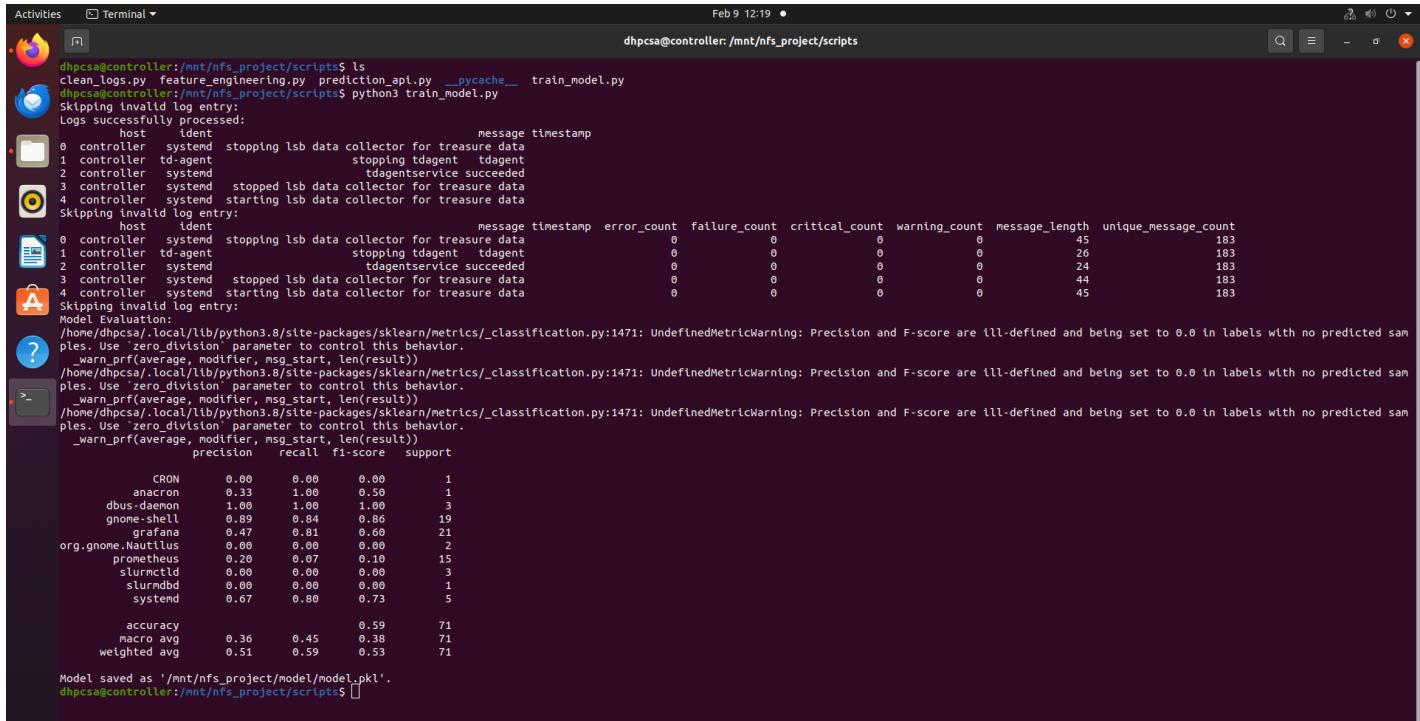
    # Train the model
    train_model(df)

if __name__ == '__main__':
    main()

dhpcsa@controller:/mnt/nfs_project/scripts$ 

```

RUNNING MODEL



```

Activities Terminal Feb 9 12:19 •
dhpcsa@controller:/mnt/nfs_project/scripts$ ls
clean_logs.py feature_engineering.py prediction_api.py __pycache__ train_model.py
dhpcsa@controller:/mnt/nfs_project/scripts$ python3 train_model.py
Skipping invalid log entry:
Logs successfully processed:
host      ident          message timestamp
0 controller system stopping lsb data collector for treasure data
1 controller td-agent      stopping tdagent  tdagent
2 controller system      tdagentservice succeeded
3 controller system stopped lsb data collector for treasure data
4 controller system starting lsb data collector for treasure data
Skipping invalid log entry:
host      ident          message timestamp error_count failure_count critical_count warning_count message_length unique_message_count
0 controller system stopping lsb data collector for treasure data 0 0 0 0 45 183
1 controller td-agent      stopping tdagent  tdagent 0 0 0 0 26 183
2 controller system      tdagentservice succeeded 0 0 0 0 24 183
3 controller system stopped lsb data collector for treasure data 0 0 0 0 44 183
4 controller system starting lsb data collector for treasure data 0 0 0 0 45 183
Skipping invalid log entry:
Model Evaluation:
/home/dhpcsa/.local/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
_warn_prfaverage, modifier, msg_start, len(result))
/home/dhpcsa/.local/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
_warn_prfaverage, modifier, msg_start, len(result))
/home/dhpcsa/.local/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1471: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
_warn_prfaverage, modifier, msg_start, len(result))
precision  recall  f1-score   support
CRON      0.00     0.00     0.00      1
anacron   0.33     1.00     0.50      1
dbus-daemon 1.00     1.00     1.00      3
gnome-shell 0.89     0.84     0.86     19
grafana   0.47     0.81     0.60     21
org.gnome.Nautilus 0.00     0.00     0.00      2
prometheus 0.20     0.07     0.10     15
slurmcld   0.00     0.00     0.00      3
slurmdbd   0.00     0.00     0.00      1
systemd    0.67     0.80     0.73      5
accuracy   0.00     0.59     0.59      71
macro avg  0.36     0.45     0.38      71
weighted avg 0.51     0.59     0.53      71
Model saved as '/mnt/nfs_project/model/model.pkl'.
dhpcsa@controller:/mnt/nfs_project/scripts$ 

```

PREDICTION SERVER SCRIPT

```
Activities Terminal ▾ dhpcsa@controller:/mnt/nfs_project/scripts$ ls
dhpcsa@controller:/mnt/nfs_project/scripts$ cat prediction_api.py
from flask import Flask, request, jsonify
import pandas as pd
import joblib
import logging
# Load the trained model (make sure this path points to the correct model file)
model = joblib.load('/mnt/nfs_project/model/log_model.pkl')

# Feature engineering function
def feature_engineering(df):
    keywords = ['error', 'failure', 'critical', 'warning'] # Keywords to look for in logs
    # Count the number of occurrences of each keyword in the 'message' column
    for keyword in keywords:
        df[keyword + '_count'] = df['message'].apply(lambda x: x.count(keyword))

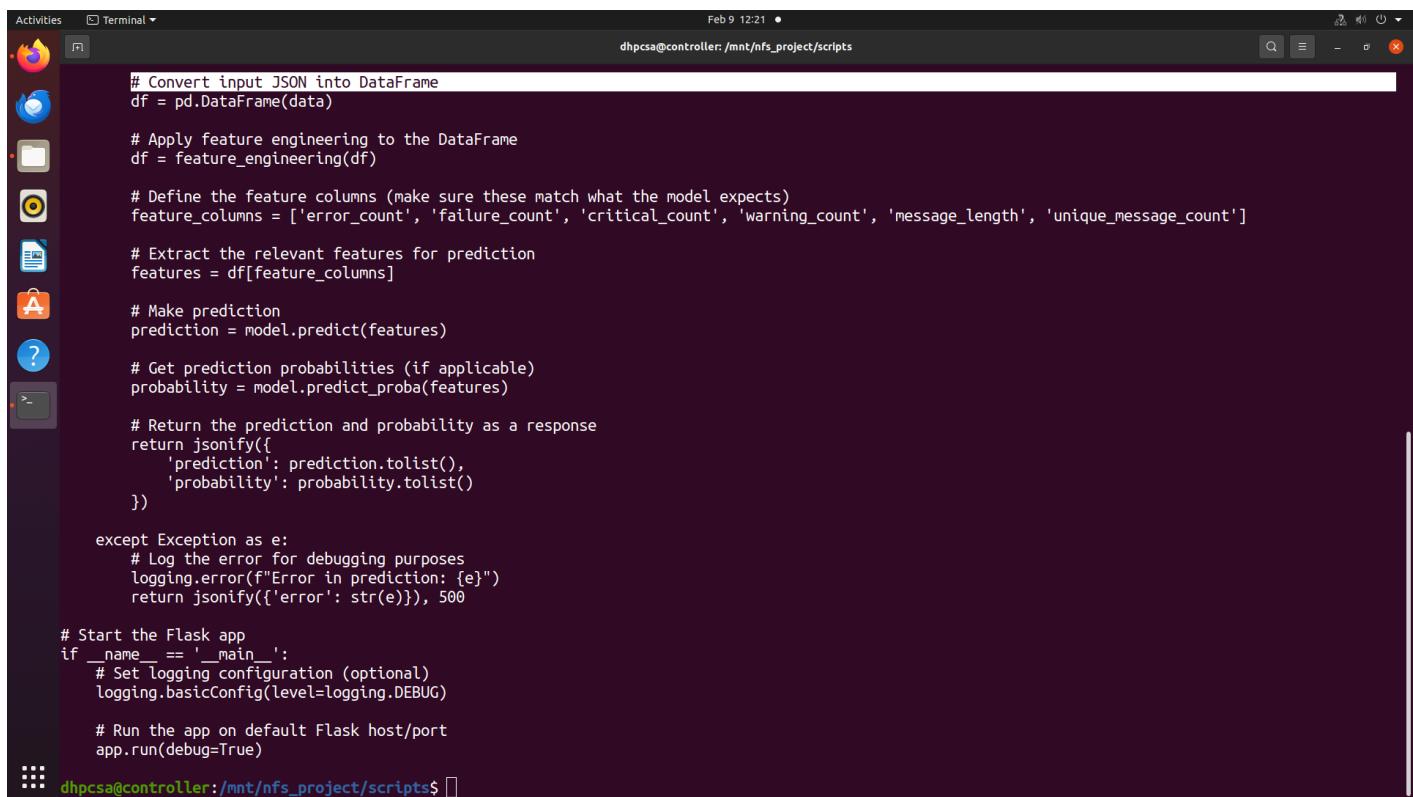
    # Add additional features
    df['message_length'] = df['message'].apply(len)
    df['unique_message_count'] = df.groupby('host')['message'].transform('nunique')

    return df

# Initialize Flask app
app = Flask(__name__)

@app.route('/predict', methods=['POST'])
def predict():
    try:
        # Get input JSON data
        data = request.get_json()

        # Check if the required 'message' field is present
        if 'message' not in data:
            return jsonify({'error': "'message' field is missing from input data"}), 400
    except Exception as e:
        # Convert input JSON into DataFrame
        return jsonify({'error': str(e)}), 500
    finally:
        # Clean up resources
        pass
```



A screenshot of a Linux desktop environment showing a terminal window. The terminal window title is "Terminal" and the path is "dhpcsa@controller:/mnt/nfs_project/scripts". The terminal content displays a Python script for a machine learning application. The script imports pandas, applies feature engineering, defines feature columns, extracts relevant features, makes predictions, and returns the results as JSON. It also handles exceptions and logs errors. Finally, it starts a Flask app if the script is run directly. The script ends with a command to run the app.

```
# Convert input JSON into DataFrame
df = pd.DataFrame(data)

# Apply feature engineering to the DataFrame
df = feature_engineering(df)

# Define the feature columns (make sure these match what the model expects)
feature_columns = ['error_count', 'failure_count', 'critical_count', 'warning_count', 'message_length', 'unique_message_count']

# Extract the relevant features for prediction
features = df[feature_columns]

# Make prediction
prediction = model.predict(features)

# Get prediction probabilities (if applicable)
probability = model.predict_proba(features)

# Return the prediction and probability as a response
return jsonify({
    'prediction': prediction.tolist(),
    'probability': probability.tolist()
})

except Exception as e:
    # Log the error for debugging purposes
    logging.error(f"Error in prediction: {e}")
    return jsonify({'error': str(e)}), 500

# Start the Flask app
if __name__ == '__main__':
    # Set logging configuration (optional)
    logging.basicConfig(level=logging.DEBUG)

    # Run the app on default Flask host/port
    app.run(debug=True)
```

FUTURE ENHANCEMENT

As the HPC Node/Cluster Failure Prediction System evolves, several enhancements can be implemented to improve its functionality, accuracy, and user experience. The following points outline potential future enhancements in detail:

1. Incorporate More Metrics

Objective: Expand the range of monitored metrics to enhance the predictive capabilities of the system.

- **Diverse Hardware Metrics:** In addition to CPU, memory, and disk I/O, consider monitoring additional hardware metrics such as:
 - **GPU Utilization:** For systems utilizing GPUs, tracking GPU load and memory usage can provide insights into potential bottlenecks or failures specific to graphical processing tasks.
 - **Network Latency and Throughput:** Monitoring network performance metrics can help identify issues related to data transfer between nodes, which may lead to job failures or slowdowns.
 - **Power Consumption:** Tracking power usage can help identify overheating issues or inefficiencies in resource utilization, which may precede hardware failures.
- **Application-Specific Metrics:** Incorporate metrics that are specific to the applications running on the HPC cluster. This could include:
 - **Job Execution Time:** Analyzing the time taken for jobs to complete can help identify anomalies that may indicate underlying issues.
 - **Error Rates:** Monitoring the frequency of application-specific errors can provide early warnings of potential failures.

- **Environmental Metrics:** Consider integrating environmental sensors that monitor conditions such as temperature and humidity in the data center. Extreme environmental conditions can lead to hardware failures, and early detection can help mitigate risks.

2. Improve Model Accuracy

Objective: Continuously refine predictive models to enhance their accuracy and reliability.

- **Incremental Learning:** Implement techniques for incremental learning, allowing the model to update itself as new data becomes available. This approach can help the model adapt to changing workloads and system behaviors over time.
- **Feature Selection and Engineering:** Regularly review and refine the features used in the predictive models. This may involve:
 - **Automated Feature Selection:** Utilize algorithms that can automatically select the most relevant features based on their predictive power, reducing noise and improving model performance.
 - **Advanced Feature Engineering:** Explore new feature engineering techniques, such as creating interaction terms or using domain-specific knowledge to derive new features that may enhance predictive capabilities.
- **Model Ensemble Techniques:** Consider using ensemble methods that combine multiple models to improve overall prediction accuracy. Techniques such as bagging, boosting, or stacking can help leverage the strengths of different algorithms.
- **Hyperparameter Optimization:** Continuously optimize model hyperparameters using techniques such as grid search, random search, or Bayesian optimization to ensure that models are performing at their best.
- **Regular Model Evaluation:** Establish a routine for evaluating model performance using a validation dataset. This should include monitoring key performance metrics (e.g., Precision, Recall, F1 Score) and adjusting the models as necessary based on performance trends.

3. User Feedback

Objective: Gather user feedback to enhance the dashboard and alerting systems, ensuring they meet the needs of administrators and users.

- **User Surveys and Interviews:** Conduct regular surveys and interviews with users to gather insights on their experiences with the dashboard and alerting systems. This feedback can help identify pain points, desired features, and areas for improvement.
- **Usability Testing:** Implement usability testing sessions where users can interact with the dashboard and provide real-time feedback on its functionality and design. Observing users as they navigate the system can reveal usability issues that may not be apparent through surveys alone.
- **Feature Requests:** Create a structured process for users to submit feature requests or suggestions for improvements. This could be facilitated through a dedicated feedback section within the dashboard or an external platform.
- **Iterative Design Improvements:** Use the feedback gathered to make iterative improvements to the dashboard and alerting systems. This may involve:
 - **Enhancing Visualizations:** Based on user preferences, improve the visualizations used in the dashboard to make them more intuitive and informative.
 - **Customizable Alerts:** Allow users to customize alert thresholds and notification preferences, ensuring that alerts are relevant and actionable.
- **Training and Documentation:** Provide ongoing training sessions and updated documentation based on user feedback. This will help users better understand the system's capabilities and how to leverage them effectively.

By focusing on these future enhancements, the HPC Node/Cluster Failure Prediction System can evolve into a more robust, accurate, and user-friendly tool that significantly improves the reliability and efficiency of high-performance computing environments.

CONCLUSION

In summary, the objectives outlined in this project provide a comprehensive roadmap for developing an HPC Node/Cluster Failure Prediction System. By focusing on proactive maintenance, log analysis, real-time monitoring, and machine learning-based prediction, the project aims to create a robust solution that addresses the critical challenges faced by HPC environments. The successful implementation of this system will not only improve the reliability of HPC clusters but also empower administrators with the tools and insights needed for proactive maintenance, ultimately leading to enhanced performance and reduced downtime.

Through rigorous testing and validation in real-world settings, this project aspires to contribute significantly to the field of high-performance computing and set a new standard for failure prediction and management. By achieving these objectives, the project aims to improve the reliability, efficiency, and maintainability of HPC clusters, ultimately reducing operational costs and ensuring uninterrupted service availability.

The integration of advanced machine learning techniques with real-time monitoring and alerting systems will provide HPC administrators with a powerful toolset for managing their resources effectively. As HPC environments continue to grow in complexity and scale, the need for sophisticated failure prediction systems becomes increasingly critical. This project not only addresses this need but also lays the groundwork for future advancements in predictive analytics within the realm of high-performance computing.

By fostering a proactive maintenance culture and leveraging data-driven insights, the project aims to transform how HPC clusters are managed, ensuring that they remain reliable and efficient in meeting the demands of modern computational workloads. Ultimately, the successful execution of these objectives will lead to a more resilient HPC infrastructure, capable of supporting the cutting-edge research and applications that drive innovation across various scientific and industrial domains.

REFERENCES

[1] High-Performance Computing

- Dongarra, J. J., & Sullivan, T. (2018). *High-Performance Computing: Modern Systems and Practices*. Morgan Kaufmann.
- Pacheco, P. S. (2011). *An Introduction to Parallel Programming*. Morgan Kaufmann.

[2] Job Scheduling and Resource Management

- Yoon, S. W., & Kim, H. J. (2015). "A Survey of Job Scheduling in High-Performance Computing." *Journal of Supercomputing*, 71(1), 1-25. DOI: 10.1007/s11227-015-1510-0.

[3] Monitoring Tools

- Prometheus. (n.d.). *Prometheus: Monitoring System & Time Series Database*. Retrieved from <https://prometheus.io/>
- Grafana Labs. (n.d.). *Grafana: The Open Platform for Analytics and Monitoring*. Retrieved from <https://grafana.com/>
- Node Exporter. (n.d.). *Prometheus Node Exporter*. Retrieved from https://github.com/prometheus/node_exporter

[4] Log Analysis and Parsing

- Fluentd. (n.d.). *Fluentd: Unified Logging Layer*. Retrieved from <https://www.fluentd.org/>
- Logstash. (n.d.). *Logstash: Collect, Parse, and Enrich Data*. Retrieved from <https://www.elastic.co/logstash>

[5] Machine Learning for Failure Prediction

- Ahmed, M., Mahmood, A. N., & Hu, J. (2016). "A Survey of Network Anomaly Detection Techniques." *Journal of Network and Computer Applications*, 60, 19-31. DOI: 10.1016/j.jnca.2015.11.016.
- Hodge, V. J., & Austin, J. (2004). "A Survey of Outlier Detection Methodologies." *Artificial Intelligence Review*, 22(2), 85-126. DOI: 10.1023/B:AIRE.0000045509.59082.9d.

[6] Data Preprocessing and Feature Engineering

- Kuhn, M., & Johnson, K. (2013). *Applied Predictive Modeling*. Springer.
- Han, J., Kamber, M., & Pei, J. (2011). *Data Mining: Concepts and Techniques*. Morgan Kaufmann.

[7] General References

- Bader, D. A., & Tharp, G. (2008). "High Performance Computing: A Practical Guide for Engineers and Scientists." *Springer*.