**Title: Large Language Model-Based Fault Detection in Transmission Lines: Comparative Analysis**

## 1. Objectives of the project

Objectives:

a. To develop a Large Language Model (LLM)-based model for the classification and detection of transmission line faults.
b. To compare the proposed model's accuracy and inference time against existing Machine Learning models.
c. To evaluate the proposed model using real-time data from power stations.

Aims of the project:

a. To provide a solution for power stations to classify the transmission line faults.
b. To ensure the proposed model achieves the best result from existing models.

## 2. Introduction

Transmission lines are critical components of modern electricity distribution networks, enabling the efficient transfer of electrical energy across vast geographical areas. They facilitate the delivery of power generated at remote plants to end users, ensuring a reliable electricity supply in various sectors of society. The design and operation of transmission lines encompass several key considerations, including conductor material, insulation levels, and voltage ratings, which optimize energy transmission while prioritizing reliability and safety. Their versatility allows them to carry both alternating current (AC) and direct current (DC), further enhancing their importance in maintaining a stable power supply. Key factors in the design of transmission lines include length, voltage, and current capacity. A careful balance among these parameters is essential to meet demand while minimizing energy losses. Additionally, the condition of transmission lines, particularly insulation levels, is crucial for preventing outages and ensuring public safety. Various types of transmission lines, such as overhead lines, underground lines, and submarine cables, each possess unique advantages and disadvantages suited to specific applications, highlighting the necessity of tailoring solutions to operational environments. This research focuses on developing an LLM-based fault detection model to improve the reliability and performance of electrical transmission lines by accurately classifying electrical faults.

## 3. Literature review

Machine Learning (ML) techniques have gained significant attention in electrical engineering, particularly for the classification and monitoring of transmission lines. Traditional methods, such as supervised learning, have been effective in various applications [1]-[3]; however, they often face challenges when handling large-scale datasets and achieving high classification accuracy. Ensemble learning techniques have emerged as a promising solution to these challenges by combining multiple ML models to enhance predictive performance [4]. Among these, the Stacking Ensemble method has shown superior accuracy and robustness in power system fault diagnosis [5], [6]. For instance, researchers have applied Stacking approaches to fault detection in transmission lines, demonstrating improved results compared to single-model methods.

Recent advancements in transmission line fault diagnosis have highlighted the integration of ML with real-time data analytics. Sun et al. [7] introduced an improved multiple SVM model optimized by a genetic algorithm, achieving an accuracy improvement of up to 11%. Their method was validated on an IEEE-30 node test system and real-world data, effectively addressing issues related to small sample sizes and generalization accuracy.

Yin et al. [8] developed a predictive decision support system using data mining techniques, correlating multi-source dynamic datasets with meteorological data to model transmission line disasters. This approach provided high-accuracy early warnings for potential failures, underscoring its relevance in regions like Bangladesh, where extreme weather poses significant risks to the power grid. Additionally, Tong et al. [9] proposed a novel transient fault detection and classification approach utilizing graph convolutional neural networks (GCN). By

incorporating spatial information from sampling sequences and topology data, their method has shown exceptional performance in real-time fault detection, offering advantages in online transmission line protection. Furthermore, Yu et al. [10] developed a diagnostic approach utilizing the Elgamal encryption algorithm, which not only improved data security but also achieved diagnostic accuracy exceeding 90%. Ma et al. [11] established a simulation model to identify both lightning and non-lightning strike faults, proposing criteria based on transient traveling wave current characteristics for intelligent fault diagnosis. Lahiri et al. [12] introduced a fault diagnosis method employing Decision Tree and Random Forest techniques, achieving 95% to 100% accuracy in identifying the location, type, and faulty phase of transmission lines across various scenarios.

Agarwal et al. [13] proposed a method for rapid fault identification in line commutated converter-based high voltage DC transmission lines, utilizing discrete Fourier transform analysis of DC current to enhance accuracy and reliability for timely trip commands to DC breakers. Additionally, the work presented by Hao et al. [14] introduced a faulted phase selection scheme that utilizes Multiscale Principal Entropy (MPE) values from fault transient voltage signals, combined with CS-SVM for high-accuracy fault detection, resilient to variations in fault location and other parameters. Lastly, the 10kV railway power transmission line simulation model developed by Yu et al. [15] employed a BP neural network to classify faults based on phase current differences, accurately determining fault types and locations.

Now-a-days, Large Language Model play a crucial role in different task especially generation and classification task. By integrating LLM-driven automatic WPT parameter optimization with an MHA-GRU network for temporal feature extraction, this framework delivers real-time, high-accuracy detection of rotor, air-gap, and stator faults in UHVDC synchronous condensers—significantly reducing false alarms and processing time compared to conventional methods [16].

Despite these developments, there remains a need for more comprehensive models that use Large Language Model to maximize their strengths and mitigate weaknesses. Our proposed model, which use Large Language Model, seeks to address these gaps by providing a more robust and accurate classification framework for transmission line data. In Bangladesh, where the electrical grid faces challenges from extreme weather conditions, hybrid models can enhance transmission line monitoring significantly [17]. By improving the reliability and efficiency of these systems, our research aims to contribute to stable power system operations, reducing economic losses, and ensuring consistent electricity supply.

## 4. Methodology

Our approach demonstrates a pre-trained large language model (LLM) backbone—specifically, an encoder adapted to six-channel time-series inputs (Ia, Ib, Ic, Va, Vb, Vc) by treating each fixed-length waveform segment as a "token" sequence. During training, these segmented embeddings are passed through a dropout layer (rate = 0.2) to mitigate overfitting and improve generalization across varying fault conditions. The regularized embedding then feed into a fully connected layer, which projects the high-dimensional Transformer output down to an intermediate feature space. A second dropout layer (rate = 0.1) further enhances robustness before the final output layer, a softmax classifier that predicts one of six labels (no-fault, LG, LL, LLG, LLL, LLLG). We optimize the entire network end-to-end using cross-entropy loss and AdamW, also use learning-rate warmup and cosine decay schedules. This pipeline ensures both the deep contextual understanding of the LLM and the lightweight efficiency of the classification head for real-time fault detection. Figure 1 shows the proposed approach of our model.

### 4.1 Dataset

The dataset used for this research is the Electrical Faults Analysis \& Classification dataset, which comprises a total of 7861 samples with ten distinct features [18]. Figure 2.1 shows the importance of each feature. The dataset includes data on different types of electrical faults, each characterized by a unique combination of outputs. For clarity, we combined the output features into a single column, representing the fault type. Table 1 summarizes the fault types and their corresponding output representations. Figure 2.2 displays the percentage distribution of various fault types present in the dataset.
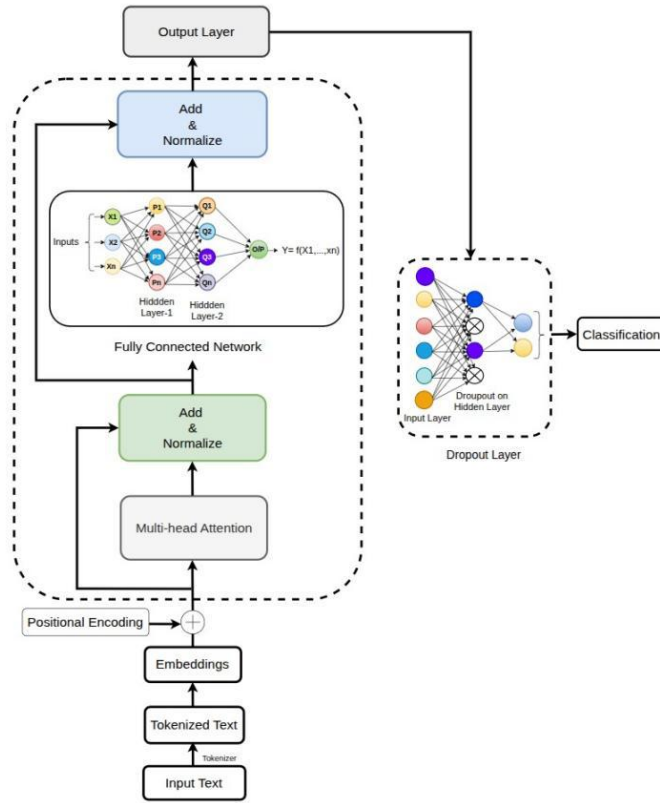
**Figure 1: Proposed Methodology**

**TABLE 1: Fault Types and their Corresponding Output Representations**

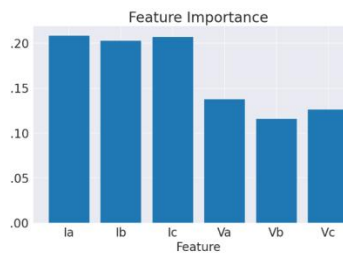| Output Representation | Fault Type |
|---|---|
| [0 0 0 0] | No-Fault |
| [1 0 0 1] | LG Fault (Between Phase A and Ground) |
| [0 0 1 1] | LL Fault (Between Phase A and Phase B) |
| [1 0 1 1] | LLG Fault (Between Phases A, B, and Ground) |
| [0 1 1 1] | LLL Fault (Between all three phases) |
| [1 1 1 1] | LLLG Fault (Three-phase symmetrical fault) |



Figure 2.1: Importance of features for classifying the fault
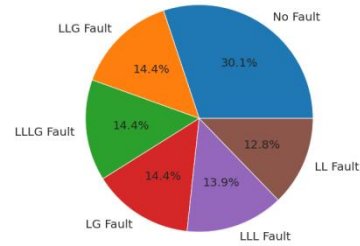


Figure 2.2: Percentage distribution of fault types in the dataset

## 4.2 Data Preprocessing

Data preprocessing is a fundamental step in preparing the dataset for ML models. In our dataset, some features have categorical data, and some feature numeric data and overfitting problems. So, in this research, we have used Label Encoding and Standard Scaler to suit the data [19], ultimately enhancing the performance and reliability of the model. Label encoding is a technique for converting categorical data into numerical format. It assigns a unique integer to each category, allowing ML algorithms to process the data. On the other hand, Standard Scaler is also a data preprocessing technique that transforms numerical features into a mean of 0 and a standard deviation of 1. This standardization process is crucial for many ML algorithms, ensuring that features with different scales contribute equally to the model's learning process.

### 4.3 LLM (Large Language Model)

We repurpose a pre-trained Transformer to process six-channel waveform data by slicing each phase's current and voltage into fixed-length segments as "tokens." Through its multi-head self-attention mechanism, the LLM captures both short-term waveform features (like fault transients) and long-range inter-phase correlations, eliminating the need for manual feature extraction.

**LLM Embedding:**

Given input tokens
$X = [x_1, ..., x_t]$ (waveform segments), the Transformer encoder produces contextual embeddings
$$H = BERT(X) \in \mathbb{R}^{T \times d}, \text{ where } H = [h_1, ..., h_t] \text{ and d is the hidden size.}$$

### 4.4 Dropout Layer

Inserted after the LLM and again after the fully connected layer, dropout randomly zeroes a set fraction of activations (commonly 10–20%) during training. This discourages co-dependent feature learning, improves robustness to noise and parameter variations, and helps the model generalize across diverse fault conditions.

Applied to embedding $h$, dropout randomly zeroes elements with probability $p$, yielding

$$\tilde{h}_i = \frac{h_i \cdot m_i}{1-p}, \quad m_i \sim Bernoulli(1-p).$$

### 4.5 Fully Connected Layer

Acting as a bridge between the Transformer outputs and classification head, this dense layer applies a learnable linear transformation followed by a nonlinearity (e.g., ReLU). It distills the LLM's high-dimensional contextual embedding into a compact feature vector that emphasizes the most discriminative patterns for fault detection.

Maps pooled embedding $h\_pool$ (e.g., the [CLS] token) to a feature vector $z$ via

$$z = \phi(W_{fc} h_{pool} + b_{fc}), \text{ where } \phi \text{ is a nonlinearity such as ReLU.}$$

### 4.6 Output Layer

A final dense layer projects the compact feature vector into a six-dimensional logits vector, one for each class (no-fault, LG, LL, LLG, LLL, LLLG). A softmax activation then converts these logits into normalized probabilities, providing clear, interpretable confidence scores for each fault type.

Produces class logits $o \in \mathbb{R}^6$ and probabilities $p$:

$$o = W_{out} z + b_{out}, \quad p_i = \frac{e^{o_i}}{\sum_{j=1}^{6} e^{o_j}}$$

### 4.7 Classification

During inference, the model selects the label with the highest softmax probability. Training uses cross-entropy loss to penalize incorrect predictions, driving the network to sharpen its decision boundaries. This end-to-end pipeline ensures rapid, accurate mapping from raw waveform segments to actionable fault categories.

The predicted fault label is

$$\hat{y} = \underset{i}{argmax} \; p_i,$$

and training minimizes the cross-entropy loss, $L$

$$L = -\sum_{i=1}^{6} y_i \log p_i$$

Where:

$L$ = Cross-entropy loss

$y_i$ = True label for sample i (typically 0 or 1 in binary classification)

$p_i$ = Predicted probability of the true class for sample i

## 5. Proposed Research Plan

The duration of this research is one year and the detail completion plan of the research is given as in Table V.

**Table V: Research Plan**

| SN | Activities/Months | 02 | 04 | 06 | 08 | 10 | 12 |
|----|-------------------|----|----|----|----|----|----|
| 1 | Literature Review | ■ | | | | | |
| 2 | Algorithms design and test | ■ | ■ | | | | |
| 3 | Evaluation and comparative Analyses | | ■ | ■ | ■ | | |
| 4 | Writing manuscript | | | ■ | ■ | ■ | |
| 5 | Writing Research Report and Submission | | | | | ■ | ■ |

## 6. Proposed Research Plan

Fault classification in electrical systems has become a crucial area of research, particularly with the growing complexity of modern power grids and the need for reliable system protection. In this study, we developed LLM-based models using BERT and other various LLM to determine the most effective approach for fault classification. Extensive experiments were conducted using electrical fault datasets to evaluate the performance of these models.

## References

[1] S. Awasthi, G. Singh, and N. Ahamad, "Classifying electrical faults in a distribution system using k-nearest neighbor (knn) model in presence of multiple distributed generators," Journal of The Institution of Engineers (India): Series B, pp. 1–14, 2024.

[2] S. Ding, M. Hao, Z. Cui, Y. Wang, J. Hang, and X. Li, "Application of multi-svm classifier and hybrid gsapso algorithm for fault diagnosis of electrical machine drive system," ISA transactions, vol. 133, pp. 529–538, 2023.

[3] R. Aiswarya, D. S. Nair, T. Rajeev, and V. Vinod, "A novel svm based adaptive scheme for accurate fault identification in microgrid," Electric Power Systems Research, vol. 221, p. 109439, 2023.

[4] M. M. Rahman, A. I. Shiplu, and Y. Watanobe, "Commentclass: A robust ensemble machine learning model for comment classification," International Journal of Computational Intelligence Systems, vol. 17, no. 1, p. 184, 2024.

[5] X. Wang and T. Han, "Transformer fault diagnosis based on stacking ensemble learning," IEEJ Transactions on Electrical and Electronic Engineering, vol. 15, no. 12, pp. 1734–1739, 2020.

[6] A. I. Shiplu, M. M. Rahman, and Y. Watanobe, "A robust ensemble machine learning model with advanced voting techniques for comment classification," in International Conference on Big Data Analytics. Springer,2023, pp. 141–159.

[7] P. Sun, X. Liu, M. Lin, J. Wang, T. Jiang, and Y. Wang, "Transmission line fault diagnosis method based on improved multiple svm model," IEEE Access, vol. 11, pp. 133 825–133 834, 2023.

[8] W. Yin, M. V. N. Gumabay, H. Lin, C. Tu, and C. Ao, "Overhead transmission lines early warning and decision support system with predictive analytics," in 2022 5th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE). IEEE, 2022, pp. 310–314.

[9] H. Tong, R. C. Qiu, D. Zhang, H. Yang, Q. Ding, and X. Shi, "Detection and classification of transmission line transient faults based on graph convolutional neural network," CSEE Journal of Power and Energy Systems, vol. 7, no. 3, pp. 456–471, 2021.

[10] H. Yu, L. Zhang, P. Zhao, Z. Liu, Z. Yang, M. Jin, and B. Hou, "Fault diagnosis of power transmission line based on elgamal encryption algorithm," in 2022 IEEE 2nd International Conference on Electronic Technology,

Communication and Information (ICETCI). IEEE, 2022, pp. 953–957.

[11] J. Ma, X. Zhang, and X. Peng, "Simulation evaluation of lightning and non lightning fault identification of transmission line," in 2022 5th International Conference on Energy, Electrical and Power Engineering (CEEPE). IEEE, 2022, pp. 452–458.

[12] S. Lahiri, C. Abhijnan, and M. A. De, "Fault diagnosis in power transmission line using decision tree and random forest classifier," in 2022 IEEE 6th International Conference on Condition Assessment Techniques in Electrical Systems (CATCON). IEEE, 2022, pp. 57–61.

[13] S. Agarwal, A. Swetapadma, C. Panigrahi, and A. Dasgupta, "Fault detection in direct current transmission lines using discrete fourier transform from single terminal current signals," in 2017 1st International Conference on Electronics, Materials Engineering and Nano-Technology (IEMENTech), 2017, pp. 1–5.

[14] F. Hao, X. Yang, G. Wang, and Y. Feng, "Transmission line fault diagnosis based on machine learning," in 2023 3rd International Conference on Consumer Electronics and Computer Engineering (ICCECE), 2023, pp. 847–850.

[15] H. Yu, M. Liu, and S. Wang, "Research on fault diagnosis in the railway power transmission line based on the modern mathematical methods," in 2018 2nd IEEE Conference on Energy Internet and Energy System Integration (EI2), 2018, pp. 1–5.

[16] A. H. Siddique, S. Tasnim, F. Shahriyar, M. Hasan, and K. Rashid, "Renewable energy sector in bangladesh: the current scenario, challenges and the role of iot in building a smart distribution grid," Energies, vol. 14, no. 16, p. 5083, 2021

[17] Zhang, D., Li, S., Hong, T., Zhang, C., & Zhao, W. (2025). Enhanced Fault Prediction for Synchronous Condensers Using LLM-Optimized Wavelet Packet Transformation. Electronics, 14(2), 308

[18] E. S. Prakash, "Electrical fault detection and classification," 2021, accessed: 2021-05-22. [Online]. Available: https://www.kaggle.com/datasets/esathyaprakash/electrical-fault-detection-and-classification

[19] B. Momotaz and T. Dohi, "Prediction interval of cumulative number of software faults using multilayer perceptron," in Applied Computing & Information Technology, R. Lee, Ed. Cham: Springer International Publishing, 2016, pp. 43–58