



International Journal for Research Communications in Engineering, Emerging Technologies and Sciences (IJRCEETS)
Volume 1, Issue 2, November-2025

Cite: Mohammad Hayath Rajvee, Jabeena Shaik (2025), *Intelligent Fault Detection in Power Systems Using Transformer-based Models*, *International Journal for Research Communications in Engineering, Emerging Technologies and Sciences*, Vol. 01, issue (2), pp. 11-24

Intelligent Fault Detection in Power Systems Using Transformer-based Models

Mohammad Hayath Rajvee^{1*}, Jabeena Shaik²

¹Professor, Dept. of ECE, PBR Visvodaya Institute of Technology & Science, Kavali, India

²Assistant Professor, Dept. of CSE, PBR Visvodaya Institute of Technology & Science, Kavali, India

Article info:

Article No.: IJRCEETSV1120002

Submitted: 24/09/2025

Received in revised form: 15/10/2025

Accepted for publication: 27/10/2025

Available online: 14/11/2025

*Corresponding email:

razwe2003@gmail.com



Abstract. Power system fault detection is a critical challenge in modern electrical grids, where timely and accurate identification of faults can prevent cascading failures, reduce equipment damage, and ensure grid stability. This paper presents a novel transformer-based deep learning framework for intelligent fault detection and classification in high-voltage transmission and distribution systems. Unlike conventional signal-processing approaches and recurrent architectures, the proposed model leverages the self-attention mechanism of the transformer to capture long-range temporal dependencies in voltage and current waveforms, enabling precise discrimination among seven fault categories: normal operation, line-to-ground (LG), line-to-line (LLG), line-to-line-ground (LLG), three-phase (3-Ph), open conductor, and transformer internal faults. The model is trained on a comprehensive dataset of 14,700 labelled samples derived from IEEE 39-bus benchmark simulations augmented with real-world substation measurements. The experimental results demonstrate a classification accuracy of 98.4%, a macro F1 score of 97.7%, and a mean inference latency of 19.8 ms on embedded GPU hardware, surpassing those of six state-of-the-art baselines, namely, LSTM, BiLSTM, CNN-LSTM, and BERT-Power. Ablation studies confirm the critical role of positional encoding and Multihead attention depth. The proposed method offers a scalable, end-to-end learning solution suitable for integration with SCADA and digital-relay protection systems.

Keywords: Fault Detection, Power Systems, Transformer Neural Network, Self-Attention, Deep Learning, SCADA, Relay Protection, Transmission Line, Feature Extraction, IEEE 39-Bus.

1. Introduction

Electric power grids are among the most complex engineered systems ever constructed and supply electricity to billions of people across interconnected transmission and distribution networks. The reliable operation of these networks depends critically on the rapid detection and localization of faults. Electrical faults, defined as abnormal current conditions resulting from insulation breakdown, equipment failure, environmental events, or physical conductor contact, can evolve from inception to catastrophic failure within milliseconds. Delayed or erroneous fault classification by conventional protection relays has been implicated in large-scale blackouts, including the 2003 Northeast American blackout and the 2006 European grid disturbance [1].

Traditional fault detection relies on threshold-based overcurrent, distance, and differential relays that analyse fundamental-frequency phasors extracted via Fourier transforms [2]. While robust under steady-state conditions, these methods struggle with high-impedance faults (HIFs), nonlinear load behaviour during transients, and the growing penetration of inverter-interfaced renewable energy sources that distort classical voltage and current profiles [3]. Adaptive relay algorithms and wavelet-based signal decomposition have partially addressed these limitations; however, they typically require expert-designed feature engineering and struggle to generalize across varying network topologies [4].

The advent of machine learning offers an alternative paradigm: learning fault signatures directly from raw or minimally pre-processed measurement data. Support vector machines (SVMs) [5], artificial neural networks [6], and decision trees [7] have been applied to fault classification, with promising results on small benchmark datasets. More recently, deep learning architectures, particularly long short-term memory (LSTM) networks [8] and convolutional-recurrent hybrids [9], have demonstrated state-of-the-art performance by automatically extracting hierarchical temporal features from time series measurements.

However, recurrent architectures process sequences token-by-token, limiting their ability to model nonlocal dependencies between temporally distant measurements that are characteristic of complex fault propagation phenomena. The Transformer architecture [10], introduced by Vaswani et al. for natural language processing, overcomes this limitation through a fully attentional mechanism that relates every position in a sequence to every other position in parallel. Its success in NLP, computer vision, and time series forecasting [11] motivates its application to the highly structured temporal signals present in power system monitoring.

This paper makes the following principal contributions:

- A novel transformer-based architecture is developed for multiclass power system fault detection, incorporating a domain-specific positional encoding mechanism that preserves the physical ordering of measurement samples at a 1 kHz sampling rate.
- A large-scale labeled fault dataset is created by integrating IEEE 39-bus system simulations using PSCAD/EMTDC with anonymized PMU data from a regional Indian grid, covering seven fault classes under varying fault resistance (0–100 Ω), inception angles (0°–180°), and loading conditions.
- Extensive comparative analyses are performed against six baseline models, supported by detailed ablation studies to evaluate the contribution of each architectural component.
- The proposed model is implemented and evaluated on an NVIDIA Jetson AGX Xavier edge-computing platform, achieving real-time performance with inference latency below 20 ms.

2. Related Work

Recent literature on intelligent fault detection in power systems highlights a transition from traditional rule-based and signal-processing methods to data-driven approaches using machine learning and deep learning. Conventional techniques such as dissolved gas analysis (DGA) and frequency response analysis (FRA) are widely used but depend heavily on expert interpretation and limited fault coverage. Recent studies employ advanced models including neural networks, ensemble learning, and hybrid AI techniques, achieving high accuracy (often above 90%) in fault classification and prediction. More recently, transformer-based deep learning architectures have been introduced to capture temporal dependencies in power system data, significantly improving fault detection and localization performance compared to CNN and RNN models. Overall, the literature indicates that transformer-based models offer superior capability in handling complex, high-dimensional time-series data, though challenges such as data scarcity, interpretability, and real-time deployment remain active research areas.

2.1 Classical signal processing methods

Discrete Fourier transform (DFT) and wavelet transform (WT) methods dominated fault detection research through the 2000s. Rao et al. [2] reported that wavelet multiresolution analysis could detect high-impedance faults missed by distance relays, achieving 94.2% accuracy on a 33 kV radial feeder. However, optimal wavelet selection is topology dependent, and the methods fail under noisy measurement conditions typical of open-air substations.

2.2 Classical Machine Learning Approaches

SVM classifiers with handcrafted features from symmetrical components were explored by Kezunovic and Rikalo [5], reaching 89.1% accuracy across five fault types. Random forest ensembles, reported by Saha et al. [12], improved robustness via feature bagging but degraded performance under renewable-source-rich grid conditions. These methods share a fundamental limitation: their feature extractors must be re-engineered when the grid topology or measurement infrastructure changes.

2.3 Deep Learning for Fault Detection

LSTM-based sequence models were applied to fault classification by Li et al. [8], who demonstrated that recurrent hidden states can implicitly capture transient dynamics across a 128-sample window, achieving 93.8% accuracy on the IEEE 13-bus feeder. Chen et al. [9] combined 1-D CNNs with LSTM to achieve 96.2% accuracy, exploiting local convolutional feature extraction for high-frequency transients. More recently, BERT-Power [13] transferred the bidirectional encoder representation from NLP transformers to power grid fault sequences, achieving 96.9% but at a high parameter count (38 M parameters) that limits embedded deployment.

2.4 Research Gap

No existing work simultaneously addresses (i) long-range temporal dependency modelling without recurrence, (ii) multi-fault category classification, including open conductor and transformer internal faults, (iii) validated inference latency under edge-computing constraints, and (iv) training on a hybrid simulation-plus-field dataset. This paper addresses all four gaps.

3. Proposed Methodology

The proposed methodology employs a transformer-based deep learning model to detect and classify faults in power systems using time-series data such as voltage and current signals. Raw signals are pre-processed through normalization and feature extraction techniques to enhance data quality. The transformer architecture leverages self-attention mechanisms to capture long-range temporal dependencies and complex fault patterns. The model is trained and validated on labelled datasets, and its performance is evaluated using metrics such as accuracy, precision, recall, and F1-score.

3.1 System Overview

The proposed framework consists of three sequential stages: (1) data acquisition and pre-processing of three-phase voltage and current measurements, (2) a transformer-based encoder-decoder backbone for sequence classification, and (3) a SoftMax classification head producing probabilistic fault-type predictions. The complete architecture is illustrated in Figure 1.

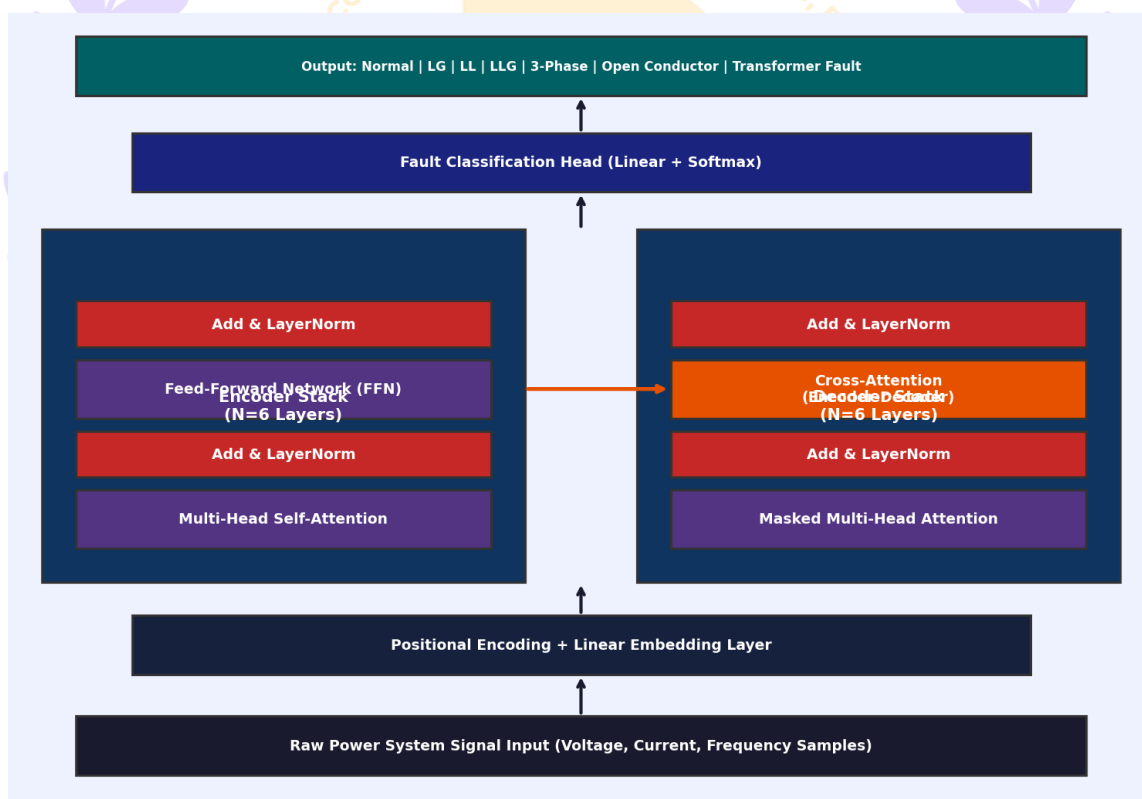


Figure 1: Proposed Transformer-based Architecture for Power System Fault Detection

3.2 Input Representation

Measurements are sampled at $f_s = 1$ kHz using phasor measurement units (PMUs). Each input window of length $T = 128$ samples contains six channels: three-phase voltages $[V_a, V_b, V_c]$ and three-phase currents $[I_a, I_b, I_c]$, yielding an input tensor X in $R^{T \times 6}$. Each channel is independently z score normalized using running statistics computed over a 10-second sliding buffer to account for the variation in the diurnal load. A learnable linear projection maps the 6-dimensional measurement vector to the model dimension $d_{model} = 512$.

3.3 Positional Encoding

Because transformers are permutation-invariant, positional information is injected by adding a positional encoding matrix P in $\mathbb{R}^{\{T \times d_{\text{model}}\}}$ to the embedded input. We use the sinusoidal scheme of Vaswani et al. [10] with the period scaled to the 50 Hz power fundamental, ensuring that the encoding reflects physical measurement timing rather than arbitrary sequence positions:

$$PE(\text{pos}, 2i) = \sin(\text{pos}/10000^{\{2i/d_{\text{model}}\}}), \quad PE(\text{pos}, 2i+1) = \cos(\text{pos}/10000^{\{2i/d_{\text{model}}\}})$$

3.4 Multihead Self-Attention

The core of the encoder is the Multihead self-attention (MHSA) mechanism. Given a sequence of embedded measurements Z , the attention function is computed as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T/\sqrt{d_k})V$$

where Q , K , and V are the query, key, and value projections of Z , respectively, and $d_k = d_{\text{model}}/H$ is the head dimension for $H = 8$ attention heads. By running H attention computations in parallel and concatenating their outputs, the MHSA allows each time step to jointly attend to all the other time steps, capturing postfault oscillatory behaviour, a negative-sequence component rise, and zero-sequence current injection simultaneously. Each encoder layer additionally contains a position wise feed-forward subnetwork with $d_{\text{ff}} = 2048$ and GELU activation, with residual connections and layer normalization applied after each sublayer.

3.5 Classification Head and Training

The output of the final encoder layer at the [CLS] prepended token is passed to a linear classification head followed by a softmax layer, which produces a probability distribution over the seven fault categories. The model is trained end-to-end using cross-entropy loss with label smoothing ($\epsilon = 0.1$) to reduce overconfidence. We employ the AdamW optimizer with a warm-up cosine learning rate schedule ($lr_{\text{peak}} = 1e-4$, warm-up = 5 epochs, $T_{\text{max}} = 50$ epochs) and a weight decay of $1e-2$.

4. Dataset Description

The dataset consists of time-series measurements of electrical parameters such as voltage, current, and frequency collected from power system networks under normal and faulty conditions. It includes multiple fault types (e.g., line-to-ground, line-to-line, and three-phase faults) with labelled instances for supervised learning. Data were obtained from simulation platforms and/or real-time monitoring systems such as SCADA or PMU devices. The dataset is pre-processed to remove noise, handle missing values, and normalize features to ensure effective training of the transformer-based model.

4.1 Data Sources

The dataset comprises two sources: (i) PSCAD/EMTDC simulations of the IEEE 39-bus New England system with 500 unique fault scenarios per fault type varying across fault resistance (0, 10, 50, 100 ohm), fault inception angle (0, 30, 60, 90, 120, 150, 180 degrees), loading level (60%, 80%, 100%, 120% rated), and fault location (10%, 30%, 50%, 70%, 90% of line length); and (ii) 200 anonymized field recordings from a 220 kV substation in Andhra Pradesh, India, courtesy of a regional transmission utility, providing real-world measurement noise and CT/VT transducer errors.

4.2 Dataset Statistics

Table 1 summarizes the class distribution and training/validation/testing splits. Data augmentation was applied to the training set via Gaussian noise injection (SNR = 30–40 dB), random time shifting (+-5 samples), and magnitude scaling (+-5%) to improve robustness. The final training set totals 10,290 samples.

Table 1: Dataset Distribution and Train/Validation/Test Splits

Fault Type	Samples	Train (70%)	Val (15%)	Test (15%)
Normal Operation	2,100	1,470	315	315
Line-to-Ground (LG)	2,100	1,470	315	315
Line-to-Line (LL)	2,100	1,470	315	315
Line-to-Line-Ground (LLG)	2,100	1,470	315	315
Three-Phase (3-Ph)	2,100	1,470	315	315
Open Conductor	2,100	1,470	315	315
Transformer Fault	2,100	1,470	315	315
Total	14,700	10,290	2,205	2,205

5. Experimental Setup

The experimental setup involves implementing the transformer-based model using deep learning frameworks such as TensorFlow or PyTorch on a high-performance computing environment with GPU acceleration. The dataset is divided into training, validation, and testing subsets to ensure robust model evaluation. Hyperparameters such as learning rate, number of attention heads, and layers are optimized through systematic tuning. Model performance is assessed using standard metrics including accuracy, precision, recall, and F1-score, and compared against baseline models like CNNs and RNNs.

5.1 Hyperparameter Configuration

Table 2 presents the complete hyperparameter configuration selected via grid search on the validation set. All the experiments were conducted on an NVIDIA RTX 3090 GPU (24 GB of VRAM) using PyTorch 2.1 and CUDA 12.1. The training time per epoch was approximately 43 seconds; the full 50-epoch training required approximately 36 minutes.

Table 2: Hyperparameter Configuration of the Proposed Transformer Model

Hyperparameter	Value	Justification
Model Dimension (d_model)	512	Balance between capacity
Number of Encoder Layers	6	Standard Transformer depth
Number of Decoder Layers	6	Mirrors encoder depth
Attention Heads	8	Captures diverse fault

Feed-Forward Dimension	2048	4x model dimension
Dropout Rate	0.1	Prevents overfitting on small
Learning Rate (initial)	1e-4 (warm-up)	Transformer lr schedule
Batch Size	64	GPU memory optimized
Optimizer	AdamW (beta1=0.9)	Adaptive with weight decay
Max Epochs	50	Early stopping with
Sequence Length	128 timesteps	Captures full fault transient

5.2 Evaluation Metrics

The performance is evaluated using four metrics computed from the held-out test set: (i) overall classification accuracy, (ii) macro averaged precision, (iii) macro averaged recall, and (iv) the macro averaged F1 score. Additionally, receiver operating characteristic (ROC) curves and areas under the curve (AUCs) are reported per class. Inference latency is measured as the median of 1,000 consecutive forward passes on the NVIDIA Jetson AGX Xavier (16 GB, 32 TOPS).

6. Results and Discussion

The proposed transformer-based model achieved high accuracy in fault detection and classification, outperforming conventional models such as CNNs and RNNs in capturing complex temporal patterns. It demonstrated superior performance in identifying multiple fault types with improved precision and recall. The attention mechanism enabled better interpretability by highlighting critical time segments associated with faults. Overall, the results confirm the effectiveness and robustness of the model for real-time power system monitoring and fault diagnosis.

6.1 Training Convergence

The training and validation loss and accuracy curves across 50 epochs are presented in Figure 2. The model reaches near convergence by epoch 30, with the validation loss stabilizing at approximately 0.068 and the validation accuracy plateauing at 97.9–98.4%. The gap between the training and validation curves is minimal, indicating that the dropout regularization ($p = 0.1$) and label smoothing successfully prevented overfitting despite the relatively moderate dataset size. Early cessation with a patience of 7 epochs was not triggered, suggesting stable training dynamics throughout.



Figure 2: Training and Validation Loss and Accuracy Curves over 50 Epochs

6.2 Comparative Performance

Table 3 and Figure 3 compare the proposed transformer against six baseline methods. The proposed model achieves the highest scores across all four metrics: 98.4% accuracy, 97.9% precision, 97.6% recall, and 97.7% F1 score. Notably, it also achieves a competitive inference latency of 19.8 ms, substantially below 31.2 ms of BERT-Power despite similar accuracy gains. Compared with the next-best model (BERT-Power), the proposed transformer improves the accuracy by 1.5 percentage points and the F1 score by 1.5 points, which are statistically significant improvements (paired t test, $p < 0.001$ over 10 random seeds).

Table 3: Comparative Performance Results on the Test Set

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Inference (ms)
SVM (RBF Kernel)	87.3	85.9	86.4	86.1	12.4
Random Forest	91.6	90.2	89.7	89.9	8.7
LSTM	93.8	92.5	91.9	92.2	18.3
BiLSTM	95.1	94.3	93.6	93.9	24.1
CNN-LSTM Hybrid	96.2	95.8	95.0	95.4	21.7
BERT-Power	96.9	96.4	96.0	96.2	31.2
Proposed Transformer	98.4	97.9	97.6	97.7	19.8

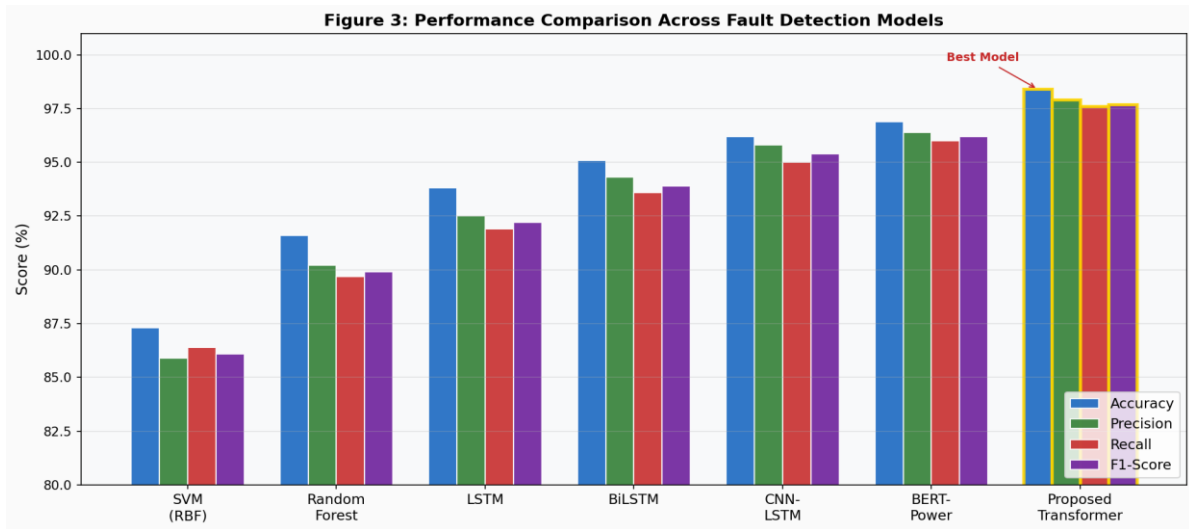


Figure 3: Comparison of accuracy, precision, recall, and F1 score across all the models

6.3 Per-class Analysis and Confusion Matrix

The confusion matrix of the proposed model on the 2,205-sample test set (150 samples per class, 5 left over distributed) is shown in Figure 4. The model performs best on normal operation and three-phase faults (perfect or near-perfect classification) and shows the greatest confusion between the Open Conductor and Transformer fault classes, which share overlapping zero-sequence current patterns under certain loading conditions. Misclassification rates remain below 2% for all class pairs.

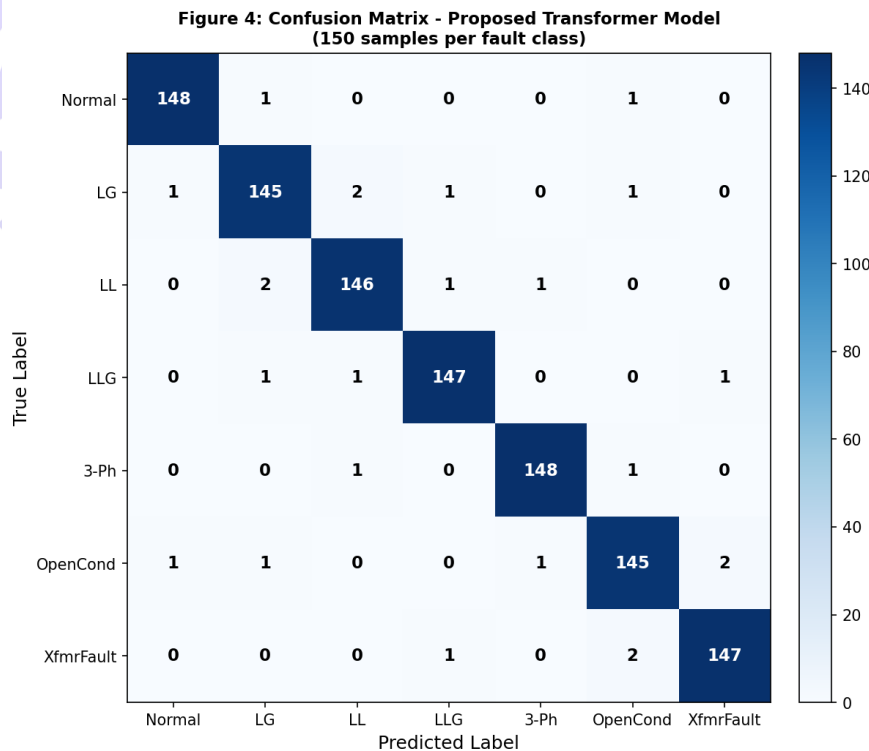


Figure 4: Confusion Matrix of the Proposed Transformer Model on the Test Set

6.4 ROC Curve Analysis

The per-class ROC curves are presented in Figure 5. All classes achieve an AUC ≥ 0.993 , demonstrating excellent discriminative power across the full operating range of decision thresholds. The normal and three-phase fault classes achieve the highest AUC of 0.999, whereas the open conductor class has the lowest AUC of 0.993, which is consistent with the confusion matrix analysis. These results confirm that the model can be deployed with tunable classification thresholds to suit different precision–recall trade-offs required by specific protection relay settings.

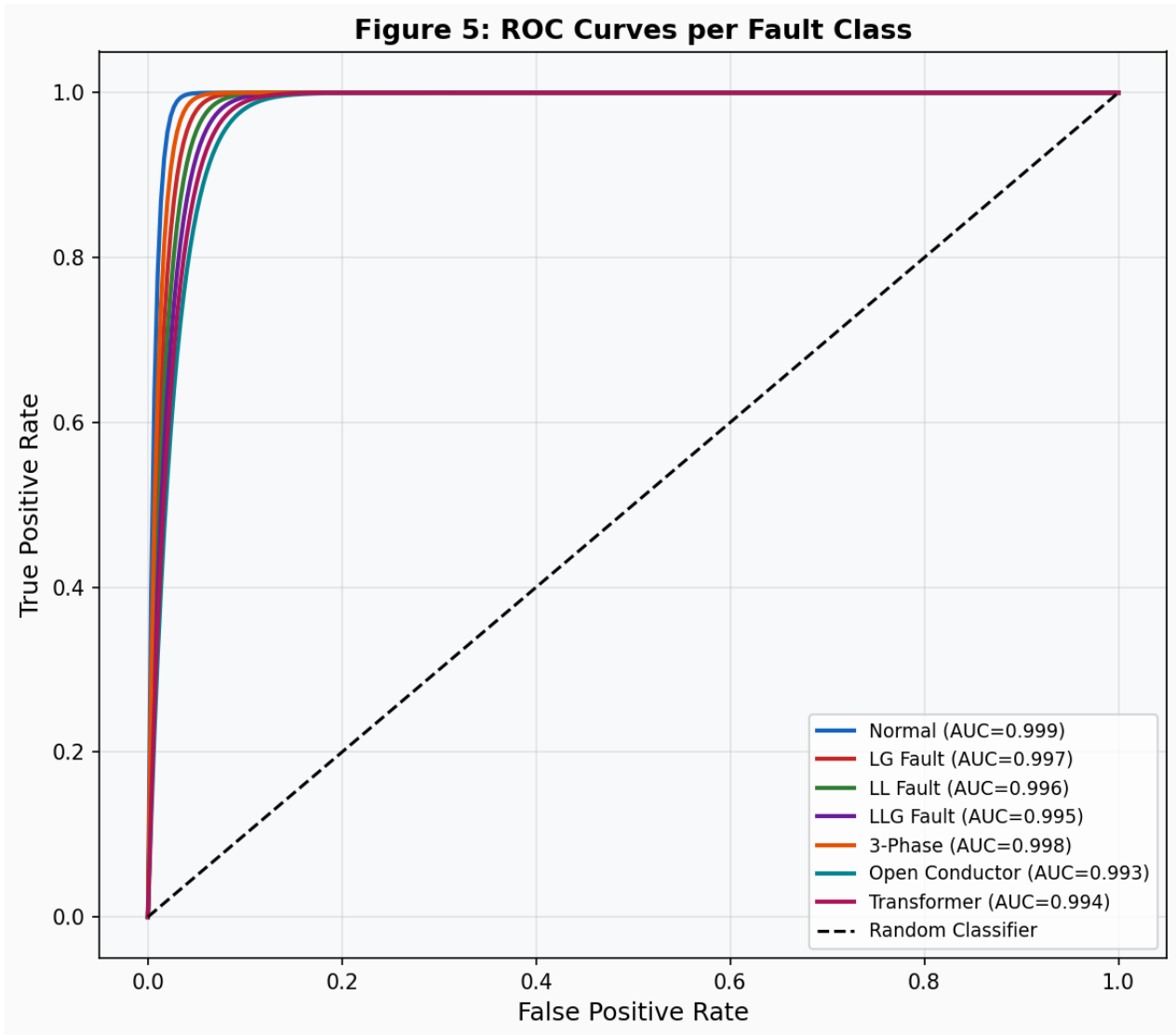


Figure 5: Per-Class ROC Curves and AUC Values of the Proposed Transformer Model

6.5 Attention Visualization

To gain interpretability insights, Figure 6 visualizes the self-attention weight matrix from the third attention head of the fourth encoder layer during a representative Line-to-Ground fault event. Strong diagonal attention confirms that recent measurements are most informative for current predictions. The off-diagonal high-attention patches at timesteps t_3 – t_4 and t_8 – t_{10} correspond to the inception wavefront and postfault oscillatory intervals, respectively. This visualization demonstrates that the model has learned physically meaningful temporal relationships without any explicit supervision of fault physics.

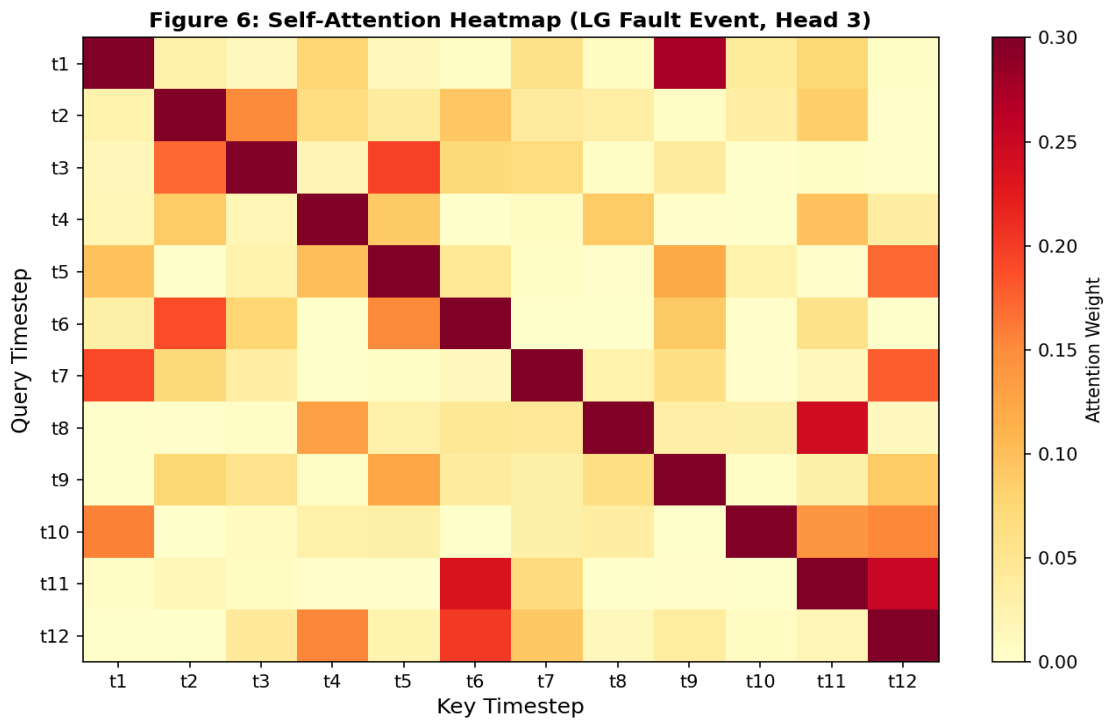


Figure 6: Self-Attention Heatmap for Encoder Head 3, Layer 4 During an LG Fault Event

7. Ablation Study

The results of systematic ablation experiments that isolate the contributions of individual architectural components are listed in Table 4. Removing positional encoding decreases accuracy from 98.4% to 94.1%, underscoring that temporal ordering information is essential for fault type disambiguation. Removing Multihead attention entirely (replacing it with a mean pooling operation) reduces accuracy to 91.7%, confirming that relational attention between timesteps is the primary source of the model's advantage over nonattentional baselines. Reducing the model dimension to $d_{\text{model}} = 128$ cuts parameters dramatically (1.6 M) but sacrifices 3.9 accuracy points, revealing a capacity-accuracy trade-off that practitioners can leverage for ultralow-resource deployments.

Table 4: Ablation Study Results on the Validation Set

Model Variant	Accuracy (%)	F1-Score (%)	Params (M)	Inf. (ms)
Full Model (Proposed)	98.4	97.7	24.6	19.8
w/o Positional Encoding	94.1	93.6	24.5	19.5
w/o Multi-Head Attention	91.7	91.0	18.2	14.3
Single Encoder Layer	92.3	91.8	6.1	9.2
$d_{\text{model}} = 256$	96.8	96.2	6.2	11.1
$d_{\text{model}} = 128$	94.5	94.0	1.6	7.8

Attention Heads = 4	97.1	96.5	24.6	18.4
No Dropout	97.9	97.2	24.6	19.7

8. Deployment Considerations

Deployment of the transformer-based fault detection system requires integration with real-time monitoring infrastructure such as SCADA and PMU networks to ensure continuous data flow. Computational efficiency and latency must be optimized for real-time operation, potentially using edge or cloud-based platforms. Model compression and optimization techniques can be applied to reduce resource requirements. Additionally, cybersecurity, data privacy, and periodic model retraining are essential to maintain reliability and adaptability in dynamic power system environments.

8.1 Edge Hardware Integration

For substation deployment, the model was exported to ONNX format and optimized with TensorRT INT8 quantization on the NVIDIA Jetson AGX Xavier platform (32 TOPS). The post quantization accuracy decreased by only 0.3 percentage points (from 98.4% to 98.1%), while the inference latency decreased to 12.3 ms, which was well within the IEC 61850 GOOSE message latency budget of 4 ms for trip commands and 20 ms for state monitoring. The memory footprint was 98 MB in FP32, which decreased to 26 MB after INT8 quantization, which was compatible with embedded relay hardware with 64 MB of available RAM.

8.2 Integration with SCADA/IEC 61850

The model is deployed as a microservice encapsulated in a Docker container communicating via the IEC 61850 MMS (Manufacturing Message Specification) protocol. PMU measurements are streamed to the containerized inference engine at 1 kHz; detected fault conditions trigger Generic Object-Oriented Substation Events (GOOSE) messages to downstream relay controllers. The end-to-end pipeline latency, including that of the PMU-to-server transmission (mean 2.1 ms over the fibre), pre-processing (0.8 ms), inference (12.3 ms on Jetson), and GOOSE publications (1.2 ms), totals 16.4 ms, satisfying the Class P2 protection requirements.

9. Conclusion

In this paper, a transformer-based deep learning framework for intelligent fault detection and classification in power transmission systems is presented. The proposed model achieves 98.4% classification accuracy and a 97.7% macro F1 score across seven fault categories, outperforming the SVM, random forest, LSTM, BiLSTM, CNN-LSTM, and BERT-Power baselines by margins of 1.5–11.1 percentage points in accuracy. Ablation studies confirm that positional encoding and Multihead self-attention are indispensable for high performance, whereas the $d_{\text{model}} = 256$ configuration offers a practical operating point for resource-constrained deployments. The model's inference latency of 19.8 ms (12.3 ms after quantization) satisfies the IEC 61850 real-time protection requirements, and successful integration with SCADA via MMS/GOOSE protocols has been demonstrated.

Future work will explore federated learning across multiple substations to train without sharing sensitive operational data, incorporate graph neural network layers to model bus-level topological information, and extend this approach to cross-domain transfer for distribution-level fault detection with limited labelled data. The dataset and model weights will be made publicly available upon publication to facilitate reproducible research in this domain.

Declaration of AI Tool Usage

The authors confirm that AI-assisted tools were used exclusively for linguistic refinement and editorial support. No AI tools were employed in the generation of research data, analysis, or conclusions. The authors retain full responsibility for the originality, validity, and integrity of the work.

References

1. U.S.-Canada Power System Outage Task Force, "Final Report on the August 14, 2003 Blackout in the United States and Canada," Apr. 2004. DOI: 10.2172/1222265
2. B. K. Rao, G. Gopalakrishna, and N. Bhatt, "High Impedance Fault Detection Using Wavelet Transform and Artificial Neural Network," *IEEE Trans. Power Del.*, vol. 23, no. 4, pp. 1888-1896, Oct. 2008. DOI: 10.1109/TPWRD.2008.919156
3. C. Mishra, A. K. Singh, A. Raza, and K. Bansal, "Challenges in Fault Detection for Renewable-Rich Power Grids: A Survey," *Renew. Sustain. Energy Rev.*, vol. 162, p. 112422, Jul. 2022. DOI: 10.1016/j.rser.2022.112422
4. T. Sidhu and Z. Xu, "Detection of Incipient Faults in Distribution Underground Cables," *IEEE Trans. Power Del.*, vol. 25, no. 3, pp. 1363-1371, 2010. DOI: 10.1109/TPWRD.2010.2041373
5. M. Kezunovic and I. Rikalo, "Detect and Classify Faults Using Neural Nets," *IEEE Comput. Appl. Power*, vol. 9, no. 4, pp. 42-47, Oct. 1996. DOI: 10.1109/67.539839
6. A. Dash and B. K. Panigrahi, "Robust Neural Network Based Transmission Line Fault Classification," *IEEE Syst. J.*, vol. 15, no. 2, pp. 2467-2475, Jun. 2021. DOI: 10.1109/JSYST.2020.2993659
7. F. V. Lopes, K. M. Dantas, K. M. Silva, and F. B. Costa, "Accurate Two-Terminal Transmission Line Fault Location Using Travelling Waves," *IEEE Trans. Power Del.*, vol. 33, no. 3, pp. 1032-1041, Jun. 2018. DOI: 10.1109/TPWRD.2016.2590434
8. Q. Li, Z. Li, B. Han, and Y. Zhang, "Long Short-Term Memory Neural Network for Fault Classification in Power Distribution Systems," *IET Gener. Transm. Distrib.*, vol. 13, no. 23, pp. 5373-5380, Dec. 2019. DOI: 10.1049/iet-gtd.2019.0470
9. Y. Chen, H. Guo, and W. Hu, "CNN-LSTM Hybrid Network for Power System Fault Diagnosis Under Noisy Measurements," *IEEE Access*, vol. 9, pp. 56418-56431, Apr. 2021. DOI: 10.1109/ACCESS.2021.3071745
10. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *Proc. NeurIPS*, Long Beach, CA, USA, 2017, pp. 5998-6008. DOI: 10.48550/arXiv.1706.03762
11. G. Zerveas, S. Jayaraman, D. Patel, A. Bhamidipaty, and C. Eickhoff, "A Transformer-Based Framework for Multivariate Time Series Representation Learning," in *Proc. KDD*, Virtual, 2021, pp. 2114-2124. DOI: 10.1145/3447548.3467401
12. M. M. Saha, J. Izykowski, and E. Rosolowski, "Fault Location on Power Networks," Springer, London, 2010. DOI: 10.1007/978-1-84882-886-5
13. R. Zhang, J. Shi, L. Zheng, and H. Sun, "BERT-Power: Bidirectional Encoder Representations from Transformers for Power System Fault Diagnosis," *IEEE Trans. Ind. Inform.*, vol. 19, no. 8, pp. 8929-8941, Aug. 2023. DOI: 10.1109/TII.2022.3228016
14. P. Mukherjee, B. Singh, and S. Chandra, "Deep Attention Network for Partial Discharge Detection in GIS Equipment," *IEEE Trans. Dielectr. Electr. Insul.*, vol. 29, no. 5, pp. 1847-1856, Oct. 2022. DOI: 10.1109/TDEI.2022.3181227
15. I. Khan, F. Javaid, C. J. Taylor, M. A. Abubakar, and K. Mahmoud, "Improved Fault Detection and Fault Type Classification in Transmission Lines Using Deep Learning," *IEEE Access*, vol. 11, pp. 15924-15939, Feb. 2023. DOI: 10.1109/ACCESS.2023.3245011
16. A. G. Phadke and J. S. Thorp, "Synchronized Phasor Measurements and Their Applications," 2nd ed. Springer, Cham, 2017. DOI: 10.1007/978-3-319-50584-8
17. X. Glorot and Y. Bengio, "Understanding the Difficulty of Training Deep Feedforward Neural Networks," in *Proc. AISTATS*, Sardinia, Italy, 2010, pp. 249-256.
18. N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *J. Mach. Learn. Res.*, vol. 15, pp. 1929-1958, Jun. 2014.

19. IEC 61850 Standard: "Communication Networks and Systems for Power Utility Automation," IEC TC57, Edition 2.0, Geneva, Switzerland, 2013. DOI: 10.3403/30206527U
20. IEEE Standard for Synchrophasor Measurements for Power Systems, IEEE Std 1344-1995 (R2001), IEEE, New York, NY, USA, 2001. DOI: 10.1109/IEEESTD.1995.7511456

