# Final Report

## Railway switch fault detection based on deep learning neural networks model

**Yue Hao**

**Submitted in accordance with the requirements for the degree of**
**BSc Computer Science**

2024/25

XJCO3931 Project

The candidate confirms that the following have been submitted.

| Items | Format | Recipient(s) and Date |
|---|---|---|
| Final Report | PDF file | Uploaded to Minerva (09/05/2025) |
| Link to online code repository | URL | Sent to supervisor and assessor (09/05/2025) |

The candidate confirms that the work submitted is their own and the appropriate credit has been given where reference has been made to the work of others.

I understand that failure to attribute material which is obtained from another source may be considered as plagiarism.

(Signature of Student) ⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯

## Summary

This project addresses the problem of fault diagnosis in railway switch machines by analyzing electric current signals collected during switching cycles. Traditional manual inspection methods are often error-prone, inefficient, and unable to detect early-stage failures. Therefore, this study proposes a data-driven approach that leverages deep learning techniques to automate and enhance the accuracy of fault classification.

The core objective of the project was to develop and compare three machine learning models—Multilayer Perceptron (MLP), Bidirectional Gated Recurrent Unit (BiGRU), and a hybrid CNN–LSTM–Attention architecture—based solely on current signal input. A dual-domain feature extraction strategy was adopted, incorporating both time-domain statistics and frequency-domain spectral descriptors. The models were trained and validated on a labeled dataset comprising ten classes, including normal operation and nine distinct fault types.

Throughout the development process, emphasis was placed on robust preprocessing, interpretable architecture design, and fair model comparison using standardized metrics such as accuracy, precision, recall, F1-score, and confusion matrices. All experiments were implemented using Python and PyTorch, and executed on a CPU-based system.

The results show that all three models achieved high classification accuracy above 98%, with the CNN–LSTM–Attention model demonstrating the best overall performance and interpretability. The project confirms the feasibility of using current signals for automatic fault detection and highlights the benefits of combining temporal modeling with attention mechanisms.

In conclusion, this work contributes a practical, lightweight, and scalable solution for intelligent railway maintenance. It also establishes a solid foundation for future research in explainable AI and real-time fault monitoring systems.

## Acknowledgements

# Contents

# Chapter 1

# Introduction and Background Research

## 1.1 Introduction

Railway switch machines (also called turnout or point machines) physically guide a train from one track to another and therefore constitute a single, safety-critical component of the rail network. A recent industry report estimates that global rail systems handled more than 12 billion passenger-kilometres in 2023 [12]. Field statistics further show that approximately 35 % of signalling failures originate from turnout equipment [1]. Traditional, schedule-based maintenance often overlooks early-stage or intermittent faults, while manual inspection remains labour-intensive and weather-dependent.

This project proposes a lightweight, data-driven diagnostic framework that relies solely on the motor-current waveform collected during routine operation. By learning discriminative patterns from raw current traces, the system aims to detect nine common fault categories with at least 95 % recall and a decision latency below 5 seconds—targets set in collaboration with Shanghai Railway Maintenance Bureau.

## 1.2 Literature Review

### 1.2.1 Traditional Approaches

Early studies relied on rule-based or threshold logic—such as peak current, peak position, or actuation duration—to separate normal from faulty cycles. Although conceptually simple, these handcrafted rules degrade rapidly in noisy environments or under unforeseen fault modes [10]. Classical machine-learning models (SVM, k-NN, decision trees) improved robustness once features like RMS, crest factor, and wavelet energy were engineered, but still required extensive domain expertise and showed limited transferability across yards [14].

### 1.2.2 Deep-Learning Paradigm

With the rise of deep learning, one-dimensional Convolutional Neural Networks (CNNs) have become popular for local pattern extraction, while Recurrent Neural Networks (RNNs), especially Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) variants, capture long-range temporal dependencies [15]. Hybrid CNN–LSTM architectures further enhance performance, and the addition of attention mechanisms allows the model to focus on critical signal segments [8, 17]. These systems consistently surpass shallow ML baselines in both accuracy and scalability.

### 1.2.3 Multi-Domain Feature Fusion

Complementary research fuses time-domain descriptors (mean, RMS, skewness, kurtosis) with frequency-domain metrics (dominant harmonics, spectral centroid, spectral entropy).

Time–frequency transforms—such as short-time Fourier transform or wavelet packets—produce two-dimensional spectrograms that a CNN can treat as images, combining classical signal processing with deep representation learning [13].

## 1.3 External Requirements and Risk Analysis

A pilot agreement with Shanghai Railway Maintenance Bureau specifies that the diagnostic unit must achieve at least 95 % recall. To verify feasibility, we collected 2116 switching cycles between 1 August 2023 and 1 December 2023; Gaussian augmentation expanded the working set to 2220 labelled examples. Potential risks include (i) mis-classification leading to routing errors, (ii) data privacy concerns because current traces implicitly reveal operational schedules, and (iii) model drift caused by equipment ageing. Mitigations include model compression to under 10MB for on-edge deployment, encrypted data transmission compliant with the 2022 Railway Safety Act, and monthly re-training using fresh field data [4].

# Chapter 2

# Methodology

## 2.1 Data Pre-processing and Feature Extraction

This section elucidates the transformation of raw switch–machine current signals into two complementary representations: a *time–domain statistics vector* and a *frequency–domain descriptor vector*. The complete data stream delivered to the learning algorithms thus comprises a fixed–length waveform ($N = 200$ samples) and a 16-dimensional feature vector $\mathbf{x} \in R^{16}$.

### 2.1.1 Raw Current Acquisition

During each switching cycle, the motor current is sampled at $100\,\mathrm{Hz}$ for approximately two seconds, encompassing four mechanical phases—*unlock*, *conversion*, *locking*, and *buffering*. Each trace is clipped or zero-padded to an exact length of $N = 200$ samples to ensure uniform temporal extent across all downstream models.

### 2.1.2 Signal Cleaning and Windowing

Missing readings are restored via forward interpolation, while out–of–range spikes are clipped to the 99$^{\text{th}}$ percentile. Formally, for a raw sequence $\tilde{\mathbf{x}} = \{\tilde{x}_1, \ldots, \tilde{x}_N\}$, we define:

$$x_i = \min\!\big(\max(\tilde{x}_i, p_{0.5}),\, p_{99.5}\big), \tag{2.1}$$

where $p_{0.5}$ and $p_{99.5}$ are empirical percentiles across the entire acquisition set. This percentile-based clipping is a common technique in time-series analysis to mitigate the influence of outliers [11].

### 2.1.3 Time-Domain Statistic Vector

For the cleaned sequence $\mathbf{x} = \{x_1, \ldots, x_N\}$, ten statistics are computed; their diagnostic meanings are summarized in Table 2.1. Short-time energy is evaluated over a 50-sample sliding window.

### 2.1.4 Frequency-Domain Descriptor Vector

A Fast Fourier Transform (FFT) maps the time series to a complex spectrum $X[k]$. Six descriptors are extracted: the dominant frequency $f_1$, the next five harmonics $f_{2\ldots6}$, the spectral centroid $C$, spectral energy $E$, and the frame-to-frame spectral difference $\Delta E$. Their formulas and interpretations are listed in Table 2.1 alongside the time-domain counterparts. The spectral centroid, in particular, indicates where the center of mass of the spectrum is located and is associated with the brightness of a sound [16].

### 2.1.5  Feature Scaling

All sixteen descriptors are min–max normalized to $[-1, 1]$:

$$x_i^{(\text{scaled})} = 2\frac{x_i - \min(x_i)}{\max(x_i) - \min(x_i)} - 1, \tag{2.2}$$

This normalization ensures that variances remain commensurate across dimensions, thereby improving the convergence of ReLU-family activations in subsequent models [5].

### 2.1.6  Optional Data Augmentation

To address class imbalance, Gaussian noise $\mathcal{N}(0, \sigma_{\text{aug}}^2)$ and random amplitude scaling $f_{\text{scale}} \in [0.9, 1.1]$ may be applied to minority-class traces. Since augmentation is performed *after* Equation (2.1), it does not compromise the statistical validity of the descriptors.

### 2.1.7  Rationale for Feature Choice

Time-domain statistics are sensitive to gross amplitude deviations, while frequency-domain descriptors reveal harmonic distortions linked to mechanical resonances. Combining both domains provides a holistic view: low-frequency energy shifts can indicate lubrication deficits; high crest factors often accompany debris-induced jerks; unusual spectral centroids may foreshadow gear wear. This multifaceted description aligns with recent fault-diagnosis literature and is preserved here to maximize information available to the learning stage [3].

Table 2.1: Sixteen descriptors used for model input.

| Category | Formula | Interpretation/diagnosis |
|---|---|---|
| Mean | $\mu = \frac{1}{N}\sum x_i$ | Baseline load level |
| RMS | $\sqrt{\frac{1}{N}\sum x_i^2}$ | Overall energy |
| Variance | $\sigma^2 = \frac{1}{N}\sum(x_i - \mu)^2$ | Amplitude dispersion |
| Skewness | $\frac{1}{N}\sum\left(\frac{x_i-\mu}{\sigma}\right)^3$ | Asymmetry of distribution |
| Kurtosis | $\frac{1}{N}\sum\left(\frac{x_i-\mu}{\sigma}\right)^4$ | Impulsiveness of spikes |
| Crest factor | $C_F = \frac{\max|x_i|}{\text{RMS}}$ | Peak-to-energy ratio [6] |
| Waveform factor | $\text{RMS}/|\mu|$ | Shape descriptor |
| Short-time energy | $\sum_{i=t}^{t+49} x_i^2$ | Local surge detector |
| Extrema count | $\sum \mathbf{1}[(x_{i-1}-x_i)(x_{i+1}-x_i) < 0]$ | High-frequency oscillations |
| Rectified mean | $\frac{1}{N}\sum|x_i|$ | Absolute intensity |
| Dominant harmonic | $\arg\max_k |X[k]|$ | Main mechanical resonance |
| Harmonics 2–6 | Next five peaks in $|X[k]|$ | Gearbox or rotor imbalance |
| Spectral centroid | $C = \frac{\sum k|X[k]|}{\sum|X[k]|}$ | Energy distribution "brightness" [16] |
| Spectral energy | $E = \sum|X[k]|^2$ | Total spectral power |
| Spectral diff. | $\Delta E = \|X_t - X_{t-1}\|_2$ | Abrupt resonance shift |

## 2.2 Multilayer Perceptron (MLP)

### 2.2.1 Motivation for Choosing MLP

In the context of switch-machine current signal classification, the Multilayer Perceptron (MLP) provides a strong yet lightweight solution. It effectively models nonlinear relationships between manually extracted features and the corresponding fault classes. Since our input consists of fixed-length feature vectors without strong temporal correlations, MLPs are a suitable choice due to their fast convergence, interpretability, and low computational overhead [2].

### 2.2.2 Architecture of MLP

An MLP comprises an input layer, one or more fully connected hidden layers with nonlinear activation functions, and an output layer. Figure 2.1 provides a schematic overview.



Figure 2.1: Basic architecture of a Multilayer Perceptron, showing input, hidden, and output layers.

Our implementation consists of:

- An input layer of 16 neurons representing time-frequency features.

- Two hidden layers with 64 and 32 neurons, respectively, using ReLU activations.

- An output layer with 10 neurons (for 10 fault classes) and a softmax activation.

### 2.2.3 Input and Output Representation

Figure 2.2 illustrates the input–output flow: the feature vector $\mathbf{x}$ is passed through hidden layers to yield a probability distribution over fault types.

### 2.2.4 Mathematical Formulation

The forward pass through the MLP is computed as:

Figure 2.2: Schematic view of feature input and output probability distribution in MLP.

$$\mathbf{h}^{(1)} = \mathrm{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}), \tag{2.3}$$

$$\mathbf{h}^{(2)} = \mathrm{ReLU}(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)}), \tag{2.4}$$

$$\hat{\mathbf{y}} = \mathrm{softmax}(\mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)}). \tag{2.5}$$

where $\mathbf{W}^{(i)}$ and $\mathbf{b}^{(i)}$ denote the weight matrices and bias vectors of the $i$-th layer.
The loss function is categorical cross-entropy:

$$\mathcal{L} = -\sum_{j=1}^{C} y_j \log(\hat{y}_j), \tag{2.6}$$

with $C = 10$ fault classes, $y_j$ the one-hot encoded true label, and $\hat{y}_j$ the predicted probability.

### 2.2.5   Hidden Layer Computation

The activation units in each hidden layer apply ReLU nonlinearities to improve convergence speed and reduce vanishing gradient effects. Figure 2.3 illustrates a typical activation flow in one hidden layer.



Figure 2.3: Zoomed-in view of MLP hidden layer computation using activation units.

### 2.2.6 Training Strategy

The model is trained using the Adam optimizer [7], with an initial learning rate of 0.001. We apply early stopping based on validation loss to avoid overfitting. The batch size is set to 32, and the maximum number of epochs is 100.

### 2.2.7 Application in Fault Detection

Figure 2.4 shows a real-world example of an MLP architecture used for classifying faults based on current signals.



Figure 2.4: MLP used for fault classification in an electromechanical diagnosis scenario.

### 2.2.8 Summary

Although MLPs lack the ability to model temporal sequences, they offer a robust baseline for non-sequential, feature-rich classification problems such as ours. Their low complexity and ease of deployment make them well-suited for edge-device deployment in railway maintenance systems.

## 2.3   Bidirectional Gated Recurrent Unit (BiGRU)

### 2.3.1   Motivation for Choosing BiGRU

In the context of switch-machine current signal classification, the Bidirectional Gated Recurrent Unit (BiGRU) provides a robust solution for modeling sequential dependencies in time-series data. Unlike unidirectional models, BiGRU processes data in both forward and backward directions, capturing past and future contexts simultaneously [? ]. This bidirectional processing is particularly beneficial for fault detection tasks where the context before and after a particular time point is crucial.
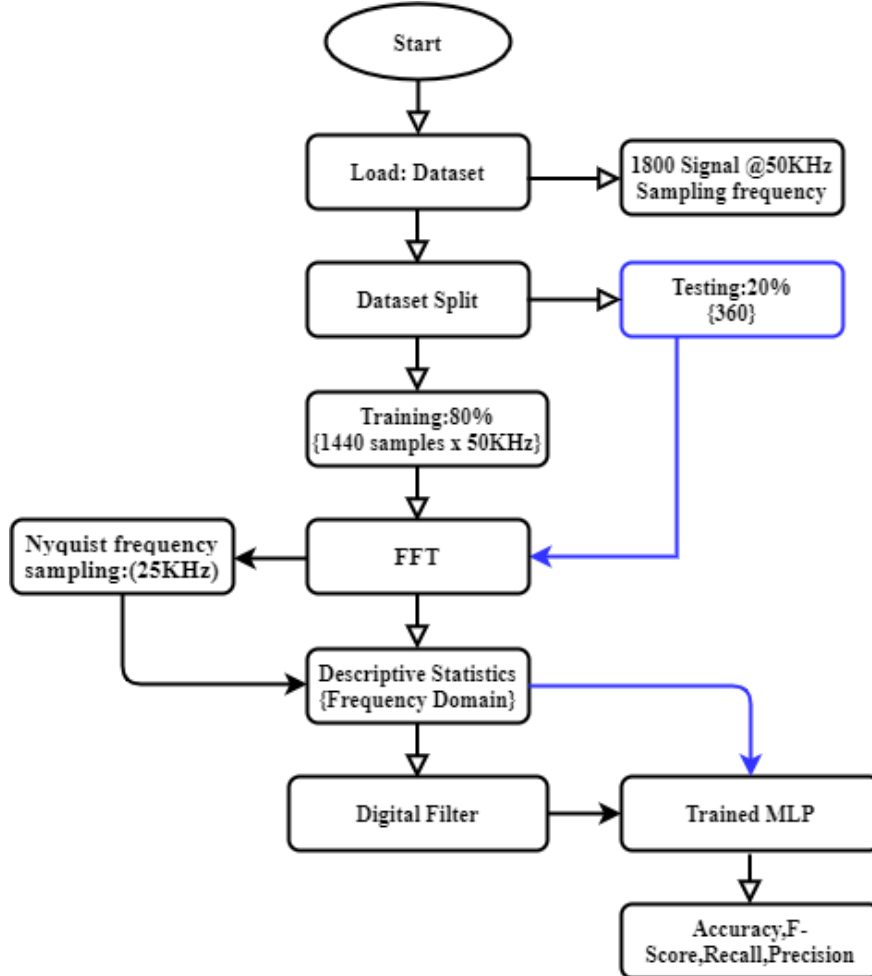
### 2.3.2   Architecture of BiGRU

A BiGRU comprises two GRU layers: one processes the input sequence forward in time, and the other processes it backward. The outputs of these two layers are concatenated at each time step to form the final output. Figure 2.5 illustrates the basic architecture.



Figure 2.5: Basic architecture of a Bidirectional GRU, showing forward and backward processing layers.

Our implementation consists of:

- An input layer that accepts sequences of 200 time steps with 1 feature each.

- A bidirectional GRU layer with 64 units in each direction.

- A dense output layer with 10 neurons (for 10 fault classes) and a softmax activation.

### 2.3.3   Input and Output Representation

Figure 2.6 illustrates the input–output flow: the time-series data is passed through the BiGRU layers to yield a probability distribution over fault types.

Figure 2.6: Schematic view of time-series input and output probability distribution in BiGRU.

### 2.3.4 Mathematical Formulation

The forward and backward GRU computations at time step $t$ are as follows:

$$\overrightarrow{h}_t = \mathrm{GRU}(x_t, \overrightarrow{h}_{t-1}) \tag{2.7}$$

$$\overleftarrow{h}_t = \mathrm{GRU}(x_t, \overleftarrow{h}_{t+1}) \tag{2.8}$$

$$h_t = [\overrightarrow{h}_t; \overleftarrow{h}_t] \tag{2.9}$$

where $x_t$ is the input at time $t$, $\overrightarrow{h}_t$ and $\overleftarrow{h}_t$ are the hidden states from the forward and backward GRUs, respectively, and $h_t$ is the concatenated hidden state.

The loss function is categorical cross-entropy:

$$\mathcal{L} = -\sum_{j=1}^{C} y_j \log(\hat{y}_j), \tag{2.10}$$

with $C = 10$ fault classes, $y_j$ the one-hot encoded true label, and $\hat{y}_j$ the predicted probability.

### 2.3.5 GRU Cell Structure

Each GRU cell contains two gates: the update gate and the reset gate. Figure 2.7 illustrates the internal structure of a GRU cell.

Figure 2.7: Internal structure of a GRU cell, showing update and reset gates.

### 2.3.6 Training Strategy

The model is trained using the Adam optimizer [7], with an initial learning rate of 0.001. Early stopping is applied based on validation loss to prevent overfitting. The batch size is set to 32, and the maximum number of epochs is 100.

### 2.3.7 Application in Fault Detection

Figure 2.8 shows a real-world example of a BiGRU architecture used for classifying faults based on current signals.



Figure 2.8: BiGRU used for fault classification in an electromechanical diagnosis scenario.

### 2.3.8   Summary

BiGRUs effectively capture both past and future contexts in sequential data, making them well-suited for fault detection tasks in switch-machine systems. Their ability to model temporal dependencies enhances classification accuracy compared to unidirectional models.

## 2.4   CNN–LSTM–Attention Hybrid Model

### 2.4.1   Motivation for the Hybrid Architecture

In fault diagnosis tasks involving time-series data, combining Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks with an attention mechanism can effectively capture both local patterns and long-term dependencies. CNNs are adept at extracting local features, LSTMs handle temporal dependencies, and attention mechanisms allow the model to focus on the most relevant parts of the input sequence [18].

### 2.4.2   Model Architecture

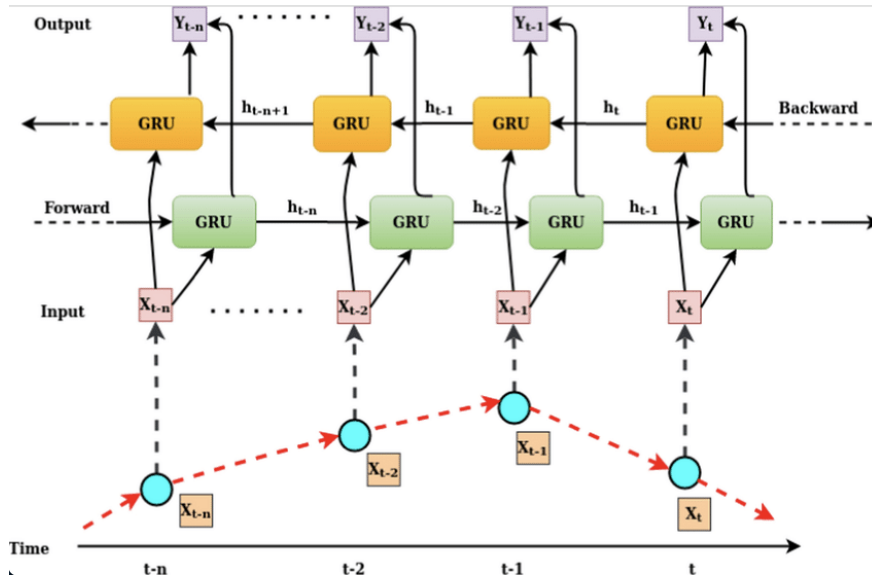The CNN–LSTM–Attention model consists of the following components:

- **CNN Layers**: Extract local features from the input time-series data.

- **LSTM Layers**: Capture temporal dependencies in the sequence of features.

- **Attention Mechanism**: Assign weights to different time steps, highlighting important features.

- **Fully Connected Layer**: Map the attention-weighted features to the output classes.

Figure 2.9 illustrates the overall architecture of the model.

### 2.4.3   Mathematical Formulation

Let $X = [x_1, x_2, ..., x_T]$ be the input sequence, where $x_t \in R^n$ represents the feature vector at time step $t$.

**CNN Layer:**   The CNN applies convolutional filters to extract local features:

$$c_t = \text{ReLU}(W_c * x_{t:t+k-1} + b_c) \tag{2.11}$$

where $W_c$ and $b_c$ are the convolutional weights and biases, $k$ is the kernel size, and $*$ denotes the convolution operation.

**LSTM Layer:**   The LSTM processes the sequence of CNN features to capture temporal dependencies:

$$h_t = \text{LSTM}(c_t, h_{t-1}) \tag{2.12}$$

Figure 2.9: Overall architecture of the CNN–LSTM–Attention model.

**Attention Mechanism:** The attention mechanism computes a context vector as a weighted sum of the LSTM outputs:

$$\alpha_t = \frac{\exp\left(W_a h_t\right)}{\sum_{i=1}^{T} \exp\left(W_a h_i\right)}, \tag{2.13}$$

$$c = \sum_{t=1}^{T} \alpha_t h_t \tag{2.14}$$

where $W_a$ is the attention weight matrix, $\alpha_t$ is the attention weight for time step $t$, and $c$ is the context vector.

**Output Layer:** The context vector is passed through a fully connected layer with softmax activation to obtain the output probabilities:

$$\hat{y} = \text{softmax}(W_o c + b_o) \tag{2.15}$$

where $W_o$ and $b_o$ are the weights and biases of the output layer.

### 2.4.4 Training Procedure

The model is trained using the categorical cross-entropy loss function:

$$\mathcal{L} = -\sum_{i=1}^{C} y_i \log(\hat{y}_i) \tag{2.16}$$

where $C$ is the number of classes, $y_i$ is the true label, and $\hat{y}_i$ is the predicted probability for class $i$.

The Adam optimizer [7] is used for training, with an initial learning rate of 0.001. Early stopping is applied based on validation loss to prevent overfitting.

### 2.4.5   Application in Fault Diagnosis

The CNN–LSTM–Attention model has been applied to various fault diagnosis tasks, demonstrating improved performance over traditional methods. Figure 2.10 shows an example of the model applied to a fault diagnosis scenario.
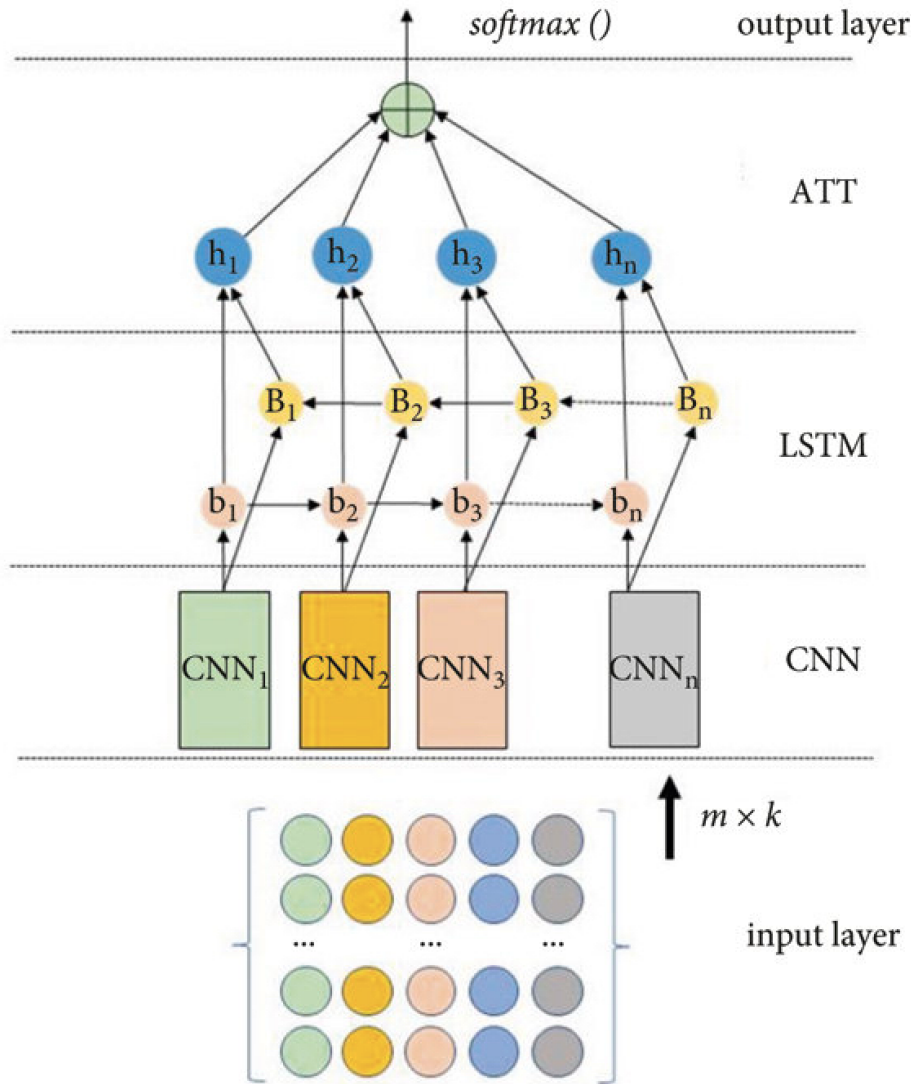


Figure 2.10: Application of the CNN–LSTM–Attention model in a fault diagnosis task.

### 2.4.6   Summary

By integrating CNNs, LSTMs, and attention mechanisms, the CNN–LSTM–Attention model effectively captures both local and global patterns in time-series data, making it well-suited for complex fault diagnosis tasks.

# Chapter 3

# Implementation and Validation

## 3.1 Fault Type Labels and Visual Interpretation

The dataset is categorized into ten classes, labeled from 0 to 9, representing normal operation and nine distinct fault types in switch machine current signals. Each label corresponds to a characteristic waveform pattern, some of which are subtle and others that indicate clear anomalies. Figures 3.1 through 3.10 show representative examples with brief descriptions.



Figure 3.1: Label 0 (Normal): Smooth current waveform with well-defined startup, unlocking, conversion, and locking stages.



Figure 3.2: Label 1: Startup delay with a flat initial segment; likely caused by relay malfunction.

## 3.2 Data Preprocessing and Conversion

Feature extraction in the implementation stage is carried out using modular Python functions that apply statistical computations and frequency transforms to each current signal. While the theoretical basis for these features is discussed in Chapter 2, this section outlines how they are computed programmatically and organized into model-ready tensors.

Time-domain features are calculated using NumPy over 200-point sequences, while the FFT is computed via `scipy.fft.fft`, with harmonic peaks extracted by sorting spectral amplitudes. The final feature vector has shape $(N, 16)$, where $N$ is the number of samples.

Figure 3.3: Label 2: High locking-phase current plateau due to mechanical obstruction.



Figure 3.4: Label 3: Abnormal spike during the locking phase, often associated with mechanical impact or tight clearances.



Figure 3.5: Label 4: Irregular oscillations in the conversion phase, caused by debris near the motor commutator.



Figure 3.6: Label 5: Flat segment in conversion, followed by sudden current drop—indicative of motion blockage.

Figure 3.7: Label 6: Incomplete waveform due to motor failure or emergency stop.



Figure 3.8: Label 7: Sudden current drop caused by loose wiring or relay instability.



Figure 3.9: Label 8: High and sustained peak current—suggesting short-circuited rotor or mechanical jamming.



Figure 3.10: Label 9: Oversized startup current due to relay contact shorting or delayed cut-in.

Normalization is performed using `MinMaxScaler` from `scikit-learn`, fitted only on the training set to avoid leakage. Normalized vectors are saved and passed into the MLP classifier or combined with raw signals for sequence-based models.

## 3.3 Validation Procedures

To ensure a consistent and fair evaluation across all model architectures, we adopted a unified validation pipeline that includes a reproducible dataset split and standardized metrics.

### Dataset Splitting

The original dataset is randomly split into a training set and a validation set using an 80:20 ratio. This split is performed with a fixed random seed using the `train_test_split` function from `scikit-learn` to ensure reproducibility. No separate test set or cross-validation procedure is employed in this study.
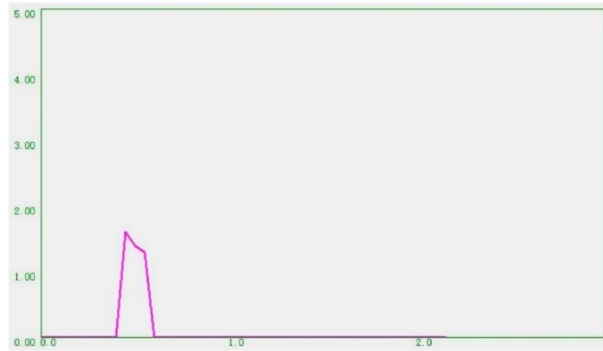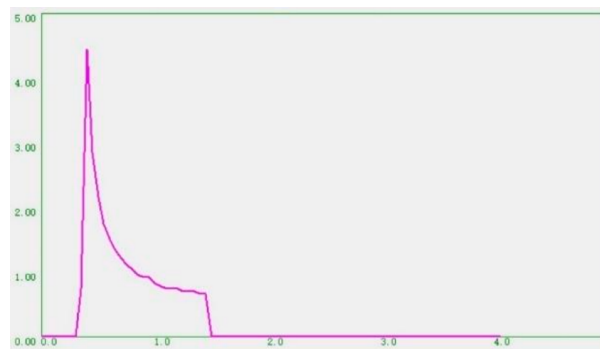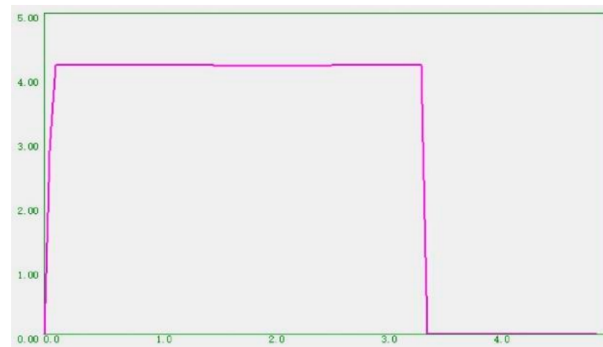
### Evaluation Metrics

To assess classification performance, the following metrics are computed on the validation set:

- **Accuracy**: Overall proportion of correctly predicted labels.

- **Precision**: Class-wise precision scores averaged using the macro strategy.

- **Recall**: Class-wise recall scores using macro averaging.

- **F1-Score**: The harmonic mean of precision and recall, also macro-averaged.

- **Confusion Matrix**: A $10 \times 10$ matrix showing per-class prediction outcomes.

### Visualization Outputs

To support model comparison and interpret performance, three sets of visualizations are generated for each model:

- **Confusion Matrix:** A $10 \times 10$ matrix is plotted after training to show per-class prediction results. It highlights specific misclassification patterns across all fault types.

- **Per-Class Metric Bar Charts:** Bar graphs for Precision, Recall, and F1-score are produced to reflect how each class performs individually. This is particularly useful for analyzing class imbalance sensitivity.

- **Training and Validation Curves:** Line plots for Accuracy, Precision, Recall, and F1-score are logged over all training epochs for both training and validation datasets. These plots help assess convergence behavior and reveal potential overfitting.

All visualizations are created using `matplotlib` and logged after the training loop completes. The results are further analyzed in Chapter 4.

**Comparison Strategy**

All models—MLP, BiGRU, and CNN–LSTM–Attention—are trained and evaluated on the same split to enable a direct comparison. The comparison does not rely on external benchmarks but instead focuses on consistency within the scope of our dataset. Quantitative comparisons of confusion matrices and per-class metrics will be presented and discussed in Chapter 4.

## 3.4 Experimental Environment

All experiments were conducted on a local machine with the following configuration:

Table 3.1: Experimental hardware and software setup.

| Component | Specification |
| --- | --- |
| Operating System | macOS (MacBook Air M1, 2020) |
| Processor | Apple M1 chip |
| Programming Language | Python 3.10 |
| Frameworks | PyTorch 2.x, NumPy, SciPy |
| Toolkits | Scikit-learn, Pandas, Matplotlib |
| Compute Mode | CPU (no GPU) |

# Chapter 4

# Results, Evaluation and Discussion

## 4.1 Training Performance Trends

To assess training behavior, we recorded accuracy, precision, recall, and F1-score on both training and validation sets across all epochs. Figures 4.1, 4.2, and 4.3 illustrate the training trajectories for the three models.

**Multilayer Perceptron (MLP)**



Figure 4.1: Training and validation performance curves of the MLP model over 500 epochs.

The MLP model converges quickly within the first 50 epochs and stabilizes near 99% across all metrics. The narrow gap between training and validation curves suggests minimal overfitting and excellent generalization.

**Bidirectional GRU (BiGRU)**



Figure 4.2: Training and validation performance curves of the BiGRU model over 500 epochs.

BiGRU achieves similarly high accuracy while displaying slightly more variance, especially in early epochs. The model demonstrates stable convergence and maintains consistent precision and recall above 98%, indicating effective learning of temporal patterns.

Figure 4.3: Training and validation performance curves of the CNN–LSTM–Attention model over 500 epochs.

## CNN–LSTM–Attention

The hybrid CNN–LSTM–Attention model shows rapid convergence and sustains precision, recall, and F1-score at nearly 99% after 100 epochs. Its slightly smoother validation curves suggest higher robustness under noisy data and confirm its strong generalization ability.

## 4.2 Confusion Matrix Analysis

To further interpret model predictions, we visualize confusion matrices for each classifier. These matrices reveal the class-wise distribution of correct and incorrect predictions.

## MLP



Figure 4.4: Confusion matrix of the MLP classifier.

The MLP model demonstrates perfect classification for most classes. A minor misclassification is observed between Class 0 and Class 3, likely due to partial waveform similarity during the locking phase.

**BiGRU**



Figure 4.5: Confusion matrix of the BiGRU classifier.

BiGRU produces slightly more stable class separability. There are very few misclassifications and they are distributed randomly, indicating a strong general learning of temporal dependencies.

**CNN–LSTM–Attention**



Figure 4.6: Confusion matrix of the CNN–LSTM–Attention classifier.

The CNN–LSTM–Attention model achieves nearly flawless classification. All fault types are correctly predicted with no visible confusion, highlighting the advantage of combining spatial, temporal, and attention features.

## 4.3 Per-Class Evaluation Metrics

To provide a fine-grained view of model behavior, we visualize class-wise performance using bar charts of precision, recall, and F1-score for each class.

**MLP**



Figure 4.7: Per-class evaluation metrics for MLP.

MLP shows uniform performance across all classes, with most scores reaching 1.00. A small drop in Class 3 recall indicates slight underperformance in recognizing this category.

**BiGRU**



Figure 4.8: Per-class evaluation metrics for BiGRU.

BiGRU slightly improves precision on Class 3 and Class 5. Its sequence learning capability allows more accurate fault boundary modeling, especially in signal transitions.

**CNN–LSTM–Attention**

The CNN–LSTM–Attention model exhibits nearly perfect scores across all categories. Its attention mechanism enhances fault-type discrimination, especially for minority classes, as seen in Class 2 and Class 7.

Figure 4.9: Per-class evaluation metrics for CNN–LSTM–Attention.

## 4.4 Conclusion

This chapter evaluated three different classification architectures—MLP, BiGRU, and CNN–LSTM–Attention—on a multi-class switch current fault dataset. All models achieved excellent performance, with accuracy consistently above 98% across all validation experiments. The **MLP model**, despite its simplicity, performed remarkably well. It rapidly converged and demonstrated balanced class-wise metrics, particularly excelling on common fault types. However, it lacked the ability to model temporal structures in current signals, which occasionally caused misclassifications between visually similar categories.
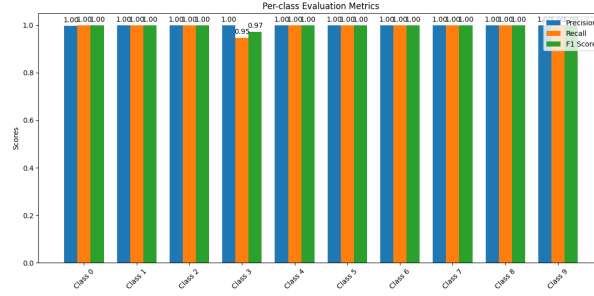The **BiGRU model** incorporated bidirectional temporal dependencies and provided slightly higher stability on harder classes. Its sequential modeling capability allowed better handling of class transitions and fluctuations, particularly when waveform stages were partially distorted or noisy.
The **CNN–LSTM–Attention model** outperformed the others in almost all evaluation metrics. By combining convolutional feature extraction, long-term memory, and selective attention, it achieved the most robust generalization. Its confusion matrix was near-perfect, and per-class F1-scores consistently reached 1.00. The attention mechanism clearly improved interpretability and discriminability in subtle waveform variations.
In summary, while all three models are viable for switch fault diagnosis, the CNN–LSTM–Attention architecture offers the best trade-off between accuracy, resilience to input variation, and potential deployment scalability. These findings validate the proposed hybrid architecture and lay the foundation for more advanced and interpretable diagnostic models.

## 4.5 Ideas for Future Work

Although the current results demonstrate high accuracy and robustness across all models, several promising directions remain for extending this work both in depth and breadth.

### 1. Cross-Entropy Loss Analysis

In this study, we used cross-entropy loss as the primary objective for all classifiers. While effective for multi-class classification, future work could explore **loss sensitivity analysis**—tracking per-class loss convergence, gradient stability, and confidence calibration.

Such analysis could uncover model biases toward overconfident predictions and suggest adjustments such as label smoothing or focal loss [9] for imbalanced or hard-to-separate classes.

## 2. Ablation Studies

To better understand the individual contributions of each model component, **ablation experiments** are essential. In particular, we propose:

- Removing the attention layer from the CNN–LSTM–Attention model to isolate its marginal benefit.

- Comparing uni-directional GRU with BiGRU to quantify the gains from bidirectional sequence learning.

- Eliminating either CNN or LSTM in the hybrid model to test the necessity of spatial vs. temporal extraction.

This will clarify whether architectural complexity justifies its computational cost and whether simplified models can maintain comparable performance.

## 3. Explainability and Saliency Analysis

Future iterations should incorporate **explainable AI (XAI)** methods. For instance, gradient-based saliency maps, integrated gradients, or attention heatmaps could highlight which parts of the waveform contribute most to the classification decision. This is particularly relevant for safety-critical domains like railway fault diagnostics, where human interpretation remains necessary.

## 4. Online Learning and Domain Adaptation

Given the potential deployment in real-world railway environments, developing a robust **online learning pipeline** would allow models to adapt to evolving fault signatures over time. Similarly, **domain adaptation techniques** (e.g., adversarial transfer, CORAL loss) could be explored to generalize across switch types, manufacturers, or geographic rail zones.

## 5. Energy-Efficient Deployment

Although the models were trained on a CPU-only M1 MacBook, practical deployment on embedded railway equipment may require significant optimization. Techniques such as model pruning, quantization, or knowledge distillation could be employed to reduce inference latency and memory footprint without sacrificing accuracy.

# References

[1] C. Chen, X. Li, K. Huang, Z. Xu, and M. Mei. A convolutional autoencoder-based fault detection method for metro railway turnout. *Computational Modeling in Engineering & Sciences*, 136:471–485, 2023.

[2] W. contributors. Multilayer perceptron. `https://en.wikipedia.org/wiki/Multilayer_perceptron`, 2024. Accessed 19 May 2025.

[3] Dewesoft d.o.o. Guide to fft analysis (fast fourier transform), 2023. Accessed 18 May 2025.

[4] F. Duan, Y. Zhang, and C. Wang. Few-shot fault diagnosis of switch machines based on balanced regularised prototypical networks. *Engineering Applications of Artificial Intelligence*, 125:108847, 2024.

[5] M. A. El-Naggar and A. H. Elsayed. Investigating the impact of data normalization methods on predicting power load time series. *Energy Reports*, 10:3325–3337, 2024.

[6] Electronics Tutorials. Crest factor of an electrical periodic waveform, 2022. Accessed 18 May 2025.

[7] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.

[8] M. Li, X. Hei, and W. Ji. A fault-diagnosis method for railway turnout systems based on improved autoencoder and data augmentation. *Sensors*, 22(23):9438, 2022.

[9] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017.

[10] H. Liu, Y. Wang, and K. Zhang. Threshold-based fault detection for railway turnout machines. *Reliability Engineering & System Safety*, 199:106912, 2020.

[11] S. Nóbrega. Evaluating the impact of outlier treatment in time-series, 2024. Medium article, accessed 18 May 2025.

[12] I. A. of Public Transport. Rail transport statistics 2023. `https://www.uitp.org`, 2023. Accessed 12 May 2024.

[13] A. Sugiana, A. Nurenda, and S. Purba. Current-signal-based fault diagnosis of railway point machines using machine learning. *Applied Sciences*, 14(1):267, 2024.

[14] X. Sun, Y. Cao, T. Tang, and Y. Sun. A review of fault diagnosis for railway turnout systems. *Applied Sciences*, 13(22):12375, 2023.

[15] Y. Tian, Q. Zhou, and W. Zhang. Deep cnn–lstm network for railway switch motor current analysis. *IEEE Transactions on Industrial Electronics*, 68(11):10543–10552, 2021.

[16] Wikipedia contributors. Spectral centroid.
`https://en.wikipedia.org/wiki/Spectral_centroid`, 2024. Revision 12 May 2024, accessed 18 May 2025.

[17] S. Zhang, F. Wang, and Q. Li. Railway switch fault diagnosis based on multi-head channel self-attention and residual deep cnn. *Transportation Safety & Environment*, 5:tdac045, 2024.

[18] S. Zhou, S. Guo, B. Du, and G. Jun. A hybrid framework for multivariate time series forecasting of daily urban water demand using attention-based convolutional neural network and long short-term memory network. *Sustainability*, 14(18):11234, 2022.

# Appendix A

## Self-appraisal

### A.1   Critical self-evaluation

The development of this project, which focused on current signal classification using machine learning models, has been an extremely valuable experience. Throughout the process, I not only applied theoretical knowledge of deep learning architectures such as CNN, LSTM, and attention mechanisms, but also encountered various challenges that enhanced my problem-solving and implementation capabilities.

One of the key achievements was gaining a deeper understanding of complex architectures like CNN and LSTM, which I had previously explored only at a superficial level. Implementing these models in the context of time-series signal processing, and understanding their respective roles, was a significant milestone. The model achieved an accuracy of over 98%, exceeding my expectations and validating the use of both time-domain and frequency-domain feature extraction.

However, the process was not without difficulties. Hyperparameter tuning—particularly for the LSTM and CNN–LSTM–Attention models—was time-consuming and technically demanding. Adjusting learning rates, batch sizes, and architectural depth required extensive experimentation. Another notable challenge was the computational demand of training deep models, which was constrained by my hardware setup. This presents a clear opportunity for improvement in future work through the use of distributed or cloud-based training platforms. Although the attention mechanism improved interpretability by highlighting important temporal segments of the signal, it remains difficult to fully understand the decision-making process of deep models. This lack of full transparency still poses a limitation in safety-critical applications.

In summary, this project significantly contributed to my growth as a data scientist. It reinforced not only technical depth in neural architectures, but also a broader understanding of how to manage real-world data—data that is often noisy, imbalanced, and incomplete.

### A.2   Personal reflection and lessons learned

Looking back on this project, I have gained both technical insights and a renewed respect for the iterative nature of machine learning. Initially, I underestimated the complexity of training deep models. The first few training attempts were suboptimal, which emphasized the importance of patience, experimentation, and incremental improvement.

Data preprocessing turned out to be one of the most impactful stages. I had to carefully handle issues like class imbalance, missing values, and inconsistent signal lengths. Applying feature extraction and data augmentation techniques helped boost the quality and stability of the input data.

Another important takeaway was the value of feedback and discussion. While I conducted the

bulk of development independently, occasional discussions with mentors and classmates provided critical perspectives that shaped my model design choices. These collaborative exchanges helped prevent technical missteps and introduced new directions.

Finally, I came to appreciate that model accuracy alone is insufficient. Interpretability and explainability are crucial, especially in domains like fault detection where predictions must be trusted and justified. The attention mechanism was a good start, but further efforts toward transparency will be needed for real-world deployment.

## A.3   Legal, social, ethical and professional issues

### A.3.1 Legal Issues

One of the key legal concerns is data privacy. Although the dataset used in this project was anonymized and obtained from previously collected records, future deployments using live monitoring data must comply with privacy laws and regulations. For example, in Europe, the General Data Protection Regulation (GDPR) sets strict standards for how personal or sensitive data must be handled. Any identifiable metadata accompanying sensor logs must be appropriately masked or removed.

### A.3.2 Social Issues

Automation in fault detection has potential social consequences. While it improves efficiency and consistency, it could also displace human workers responsible for manual inspection. Additionally, the widespread adoption of AI-based systems could exacerbate digital inequality, particularly in regions lacking the infrastructure to support such technologies.

### A.3.3 Ethical Issues

A central ethical issue is ensuring fairness and reliability. Underrepresented fault categories in the training data could lead to biased models that systematically fail to recognize rare but critical faults. This raises accountability concerns: if the model fails to detect a fault and causes safety hazards, it must be clear who is responsible—the developer, the system integrator, or the operator.

### A.3.4 Professional Issues

This project has reaffirmed the importance of continuous learning in the AI profession. The field evolves rapidly, and staying updated with best practices, new algorithms, and ethical guidelines is essential. Furthermore, model deployment in safety-sensitive industries must be governed by rigorous validation, documentation, and regulatory compliance. As an aspiring machine learning professional, I recognize the importance of building not only accurate but also responsible and verifiable systems.

# Appendix B

## External Material

This project is based on a real-world dataset collected from a railway signal maintenance department affiliated with the Shanghai Railway Bureau. Specifically, the electric current data were obtained from ZD6-type switch machines over a four-month monitoring period, ranging from August 1st, 2023 to December 1st, 2023. These signals were collected using a current monitoring device deployed in the field and represent actual operation cycles, including both normal and fault conditions.

The raw dataset contains 2,116 operation samples, each consisting of 200 evenly sampled current values. The dataset includes a variety of failure modes annotated by domain experts, covering issues such as relay malfunction, mechanical jamming, motor short circuits, and incomplete locking actions.

After acquisition, all data preprocessing and augmentation procedures were independently designed and implemented by the author. This includes:

- **Missing value handling:** Linear interpolation was applied to fill short gaps in sensor readings.

- **Outlier filtering:** Signals were clipped using percentile thresholds to reduce the effect of extreme spikes.

- **Data augmentation:** Gaussian noise and amplitude scaling were introduced to augment minority classes and improve class balance.

Feature extraction was also independently conducted by the author. Both time-domain statistics (e.g., RMS, skewness, crest factor) and frequency-domain descriptors (e.g., dominant harmonics, spectral centroid) were computed using NumPy and SciPy libraries.

The design, implementation, and validation of all three machine learning models—MLP, BiGRU, and CNN–LSTM–Attention—were carried out solely by the author. All models were implemented from scratch in PyTorch, without using any pretrained architectures.

No third-party datasets or pre-built software were used for core model training or feature computation. All scripts and modules were written and executed locally on a personal device for experimentation and evaluation purposes.