

A FAULT DIAGNOSIS METHOD OF ACTIVE BRAKING SYSTEMS BASED ON 1D-CNN AND MHAM

Hui Du

School of Computer and Artificial Intelligence, Beijing Technology and Business University, Beijing 102488, Beijing, China.

Abstract: To overcome insufficient feature representation and limited status quantification in active braking systems (ABS), an end-to-end fault diagnosis method of ABS is constructed to combine a one-dimensional convolutional neural network (1D-CNN) and a multi-head attention mechanism (MHAM) under conditions of stronger noise interference and multi-source nonlinear coupling. Based on a 1D-CNN front-end encoder to directly decouple local high-frequency transient features from the original time-series waveform, the long-distance dependency topology within the sequence is deeply reconstructed to quantify the contribution weights of features at different time steps to the evolution of the system's health status. according to the joint constraints of the Adam adaptive optimization algorithm and regularization penalty term, the model effectively avoids the overfitting risk of deep networks and significantly enhances its generalization robustness under unknown and complex conditions. Empirical results strongly demonstrate the superiority of this method. The regression prediction determination coefficient R^2 for the degree of failure reaches 0.9472, and the root mean square error (RMSE) is reduced by 1.56%, achieving high-precision quantitative perception of the system's health status.

Keywords: MHAM; Fault diagnosis; Active braking system; 1D-CNN

1 INTRODUCTION

With the rapid development of the automotive industry, the active braking system (ABS) has become the core system for ensuring driving safety. As a multi-physics coupling system of telecommunications, the ABS will be affected by various factors, such as friction during long-term operation. The serious nonlinear and time-varying decay characteristics of ABS can cause obstacles in the braking process of the vehicle and thus cause serious accidents [1]. Therefore, it is of great significance to propose a diagnostic method that can perceive the health status of the ABS in real-time and accurately quantify the severity of the faults to enhance the safety of the driving system [2].

In the field of fault diagnosis, data-driven diagnostic methods have received significant attention because the above methods do not require precise physical models. In recent years, deep learning algorithms, represented by 1D-CNN, have been widely used in fault identification in complex industrial processes [3]. Especially for one-dimensional time series signals such as vibration and pressure, 1D-CNN can directly extract features from the original waveform [4], avoiding the reliance on expert experience in manual feature extraction, and has a great advantage in processing sensor sequence data.

However, despite the excellent performance of 1D-CNN in extracting local features, it still faces a huge challenge in dealing with the long-distance temporal dependencies of multi-source signals in ABS. Since the operation of convolution is limited by a fixed receptive field, it is difficult to capture nonlinear coupling features with a large time span [5]. Simply expanding the receptive field by stacking convolutional layers will not only cause gradient vanishing and network degradation, but also increase the computational complexity, limiting the robustness and generalization of the model in complex conditions. To address the above limitations, the attention mechanism provides a novel solution. This mechanism is proposed by Vaswani et al. in the field of machine translation. It obtains different weights at different positions in the sequence through parallel computation, thereby capturing global dependencies [6]. Meanwhile, some scholars have attempted to introduce the attention mechanism into industrial fault diagnosis, proving its effectiveness in suppressing noise interference and enhancing fault features [7].

Based on this, a fault diagnosis method is proposed to combine the one-dimensional convolution and the multi-head attention mechanism (1D-CNN-MHAM). First, 1D-CNN is used to extract high-frequency features of multi-source signals. Second, the MHAM is used to obtain global temporal correlation and to achieve an accurate quantitative diagnosis of the fault degree of the ABS.

2 FAULT CHARACTERISTICS AND CLASSIFICATION OF ABS

2.1 Fault Classification of ABS

The deep electromechanical-hydraulic coupling inherent in ABS determines that their fault evolution mechanisms exhibit significant multi-source heterogeneity. Within the research context of engineering fault-tolerant control and precision diagnostics, a classification paradigm based on a control hierarchy architecture has been established as the core analytical framework due to its superior physical interpretability. This paradigm not only achieves precise

decoupling of the intricate fault propagation topology but also lays a solid logical foundation for constructing high-fidelity, high-precision fault diagnosis models.

Perception layer failures primarily stem from sensor performance degradation or external interference. For example, radar sensor miscalibration can occur due to a change in the mounting bracket's position after a front bumper accident repair, leading to a limitation of system functionality. Decision layer failures originate from hardware, algorithm, or system logic defects in the control unit. Execution layer failures involve the malfunction of mechanical, hydraulic, or electronic control components within the braking system itself, such as brake disc deformation, causing steering wheel vibration and reduced braking efficiency during braking. Table 1 shows the main manifestations of these failures and their corresponding functional defects.

Table 1 Fault Manifestations and Corresponding Functional Defects

Fault Classification	Fault symptoms	Functional defects
Perception layer fault	Environmental interference sensor dirt calibration inaccuracies Control logic defects	False alarms and missed detections The functionality has been limited Decision error.
Decision-making failure	interaction conflicts hardware and software malfunctions	Unexpected function exit Incorrect command.
Execution layer failure	Braking system malfunction, abnormal component damage abnormal coordination function	Insufficient braking force Increased braking distance Loss of vehicle stability

2.2 Fault Feature Extraction Method for ABS

The fault feature extraction stage of an active braking system involves identifying and extracting key sensitive indicators for different fault models from system operation data, and is a crucial step in constructing a fault diagnosis model. This section will describe the fault feature extraction method to be used in this study from three aspects: data extraction type, preprocessing method, and feature extraction algorithm.

2.2.1 Fault characteristic data types

The fault characteristic data of the active braking system mainly come from various sensors and system operating status signals, including high-frequency acceleration in the vibration signals of the actuator layer, such as brake disc deformation or component wear. The operating status of the actuator directly reflects the establishment of brake hydraulic pressure. By analyzing the difference between the commands issued by the electronic control unit (ECU) and the actual feedback, faults in the decision-making layer and the execution layer can be effectively diagnosed [8]. This paper mainly conducts simulation experiments based on these characteristic data.

2.2.2 Data preprocessing

Considering the significant differences in physical dimensions and orders of magnitude between signals from multiple heterogeneous sensors, directly feeding them into a neural network would inevitably induce severe oscillations in the gradient update process and even convergence lag. Therefore, this study employs a Min-Max normalization strategy, forcing all feature vectors to be linearly mapped to the closed interval [0,1]. This paradigm effectively smooths the optimization surface of the loss function while completely eliminating dimensional barriers, thus significantly accelerating the model's training and convergence process. The normalization calculation formula is shown in (1).

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (1)$$

where X_{min} and X_{max} are the minimum and maximum values of this feature dimension.

The Z-score normalization calculation is expressed by (2).

$$X_{std} = \frac{X - \mu}{\sigma} \quad (2)$$

2.2.3 Hierarchical feature extraction algorithm

Abandoning the traditional paradigm of fragmented signal correlations through multi-branch structures, this study delves into the strong nonlinear coupling mechanism between sensor channels, thereby constructing a serial fusion feature extraction architecture. In this architecture, the 1D-CNN is deployed as a primary feature extraction operator, directly performing end-to-end temporal modeling of the multidimensional joint signal.

At the local feature encoding level, given the dimensionality redundancy and computational inefficiency of two-dimensional convolution when processing temporal topology, this study establishes the methodological advantage of 1D-CNN in capturing local fluctuation patterns in fixed-length sensor sequences. Its core lies in utilizing the sliding operation of the convolution kernel in the time dimension to accurately characterize the transient features of the signal. For the l -th layer of convolution, its operation is defined by (3).

$$y_k = \sum_{i=0}^{K-1} w_i \cdot x_{t+i} + b \quad (3)$$

where w is the convolution kernel weight, and b is the bias.

To alleviate the gradient vanishing problem and accelerate convergence in deep networks, this study introduces a Batch Normalization (BN) layer after the convolution operation. This layer normalizes the features to stabilize their distribution, and its calculation process is shown by (4).

$$\hat{y}_k = \frac{y_k - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, z_k = \gamma \hat{y}_k + \beta \quad (4)$$

where μ_B and σ_B^2 are the mean and variance of the current batch of data, and γ and β are learnable scaling and translation parameters. Subsequently, a nonlinear mapping is introduced through the ReLU activation function to enhance the model's ability to fit complex working conditions.

Global Dependency Mining: Although the local feature sequences extracted through convolution operations are rich in short-term dynamic information, they suffer from a natural bottleneck in constructing long-span temporal dependencies due to the locality of the convolution kernel's receptive field. Therefore, this study cascades a multi-head self-attention mechanism at the back end of the convolutional feature extractor to overcome the limitations of a local perspective. This mechanism adaptively reconstructs the topology of association weights between time steps within the sequence by mapping the high-dimensional feature space to three latent subspaces: query, key, and value. This allows for the deep analysis of implicit long-distance coupled fault modes from a global perspective. The mathematical expression of its core attention score is obtained by (5).

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} V \right) \quad (5)$$

where d_k is the dimension of the key vector, and $\sqrt{d_k}$ is used to scale the dot product result to prevent gradient vanishing.

Finally, to reduce feature dimensionality and extract the most significant fault features, a Global Average Pooling (GAP) layer is used to aggregate the features along the time dimension. For a feature sequence H of length L , the GAP operation is defined by (6).

$$v = \frac{1}{L} \sum_{i=1}^L H_i \quad (6)$$

The generated compact feature vector v is then fed into a fully connected layer for regression prediction of the degree of failure.

3 FAULT DIAGNOSIS METHOD FOR ABS BASED ON MHAM

3.1 Basic Modeling Methods of MHAM

As an advanced modeling technique that can deconstruct the multi-subspace dependencies of sequences in high-dimensional space through parallel computing, the MHAM is strategically deployed after the feature extraction layer of 1D-CNN in this paper. This design aims to deeply explore the long-distance nonlinear coupling mechanism between active braking system signals implied in the temporal feature sequence after convolutional encoding, and dialectically quantify the difference in contribution weight of features at different time steps to the evolution of the overall health state of the system [9].

3.1.1 Core principles and structure

The theoretical core of the multi-head attention mechanism is based on the high-order generalization of the scaled dot-product self-attention operation. Unlike the limitations of single-channel attention, this mechanism forces the input feature sequence to be orthogonally projected to a multi-dimensional feature subspace through multiple independent linear transformation matrices. This design gives the model the dialectical ability to analyze heterogeneous information in parallel in different dimensions. By high-dimensional splicing and quadratic linear projection of the output vectors of all subspaces, the model can reconstruct a complete feature representation that has both local fineness and global macroscopic vision [6].

3.1.2 Mathematical model and solution parameters

Let the feature sequence obtained by the preceding 1D-CNN module be $H \in \mathbb{R}^{L \times d_{model}}$ where L is the sequence length and d_{model} is the feature dimension. The operation logic of the multi-head attention layer follows these steps.

Generate independent Q, K, V matrices for each head i , which is achieved through learnable linear projection by (7).

$$Q_i = HW_i^Q, K_i = HW_i^K, V_i = HW_i^V \quad (7)$$

where $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{model} \times d_k}$ is the projection weight matrix specific to the i -th head. Typically, $d_k = d_{model}/h$ is set to balance computational cost.

Then, the scaled dot product attention is computed in parallel within each head by (8).

$$\text{head}_i = \text{Attention}(Q_i, K_i, V_i) = \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_k}} \right) V_i \quad (8)$$

Building upon this, the scaled dot product attention is computed in parallel within each subspace. To avoid the Softmax function falling into the saturation region (Vanishing Gradients) due to excessively large dot product results, a scaling factor $\sqrt{d_k}$ is introduced to normalize the correlation matrix, thereby ensuring the numerical stability of the backpropagation process. Following this, the outputs of all h heads are concatenated along the feature dimension by (9).

$$\text{MultiHead}(H) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) \quad (9)$$

By performing linear projection using the second key learnable parameter matrix W^O , the feature dimensions are recovered and the information from each head is fused to obtain the final global feature representation by (10).

$$O = \text{MultiHead}(H)W^O \quad (10)$$

3.1.3 Parameter setting and solution model

The core trainable parameter set of this module is defined as $\theta_{\text{att}} = \{W_i^Q, W_i^K, W_i^V, W_i^O\}_{i=1}^h$. In the overall training process of the model, the above parameters are tightly coupled with the front-end CNN feature extraction network and the back-end fully connected regression layer, and end-to-end joint learning is carried out by relying on the backpropagation algorithm and the Adam adaptive optimizer [10]. Through the deployment of this mechanism, the features of each time step in the output sequence are deeply integrated with the dynamic weighted information from the global context of the original multi-source sequence, thereby constructing a high-order description of the system fault state and providing complete information for subsequent accurate diagnosis.

3.2 Fault Diagnosis Method Based on MHAM

This section integrates the MAHA into a complete end-to-end deep learning framework, aiming to solve the fault severity diagnosis problem in active braking systems. Unlike simple fault classification, this method can output continuous health indicators (such as remaining brake pad thickness and braking efficiency coefficient), providing a more accurate basis for condition-based maintenance (CBM) decisions.

3.2.1 Methodological Framework and Solution Process

This method adopts an end-to-end supervised learning approach. Its core process includes extracting local temporal features from multi-source raw signals using 1D-CNN [11], performing global feature fusion and long-distance dependency mining through the MHA layer, and obtaining the final fault degree prediction value through global average pooling and the fully connected regression layer mapping. The specific solution flowchart is shown in Figure 1.

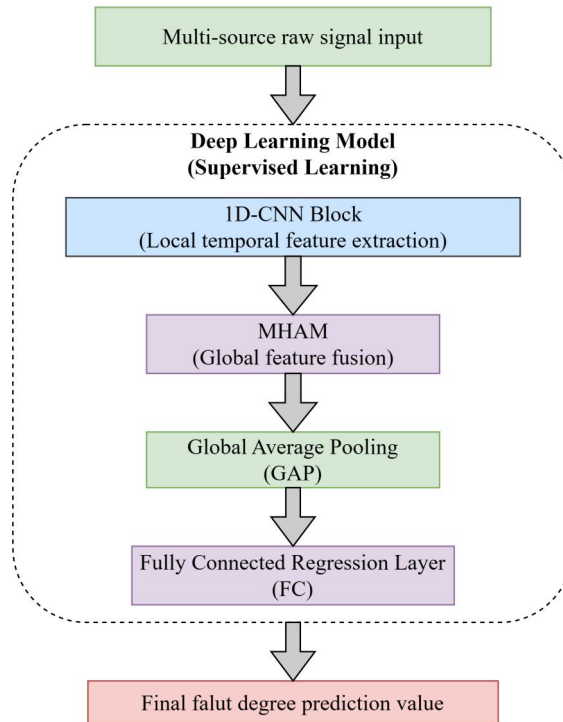


Figure 1 The Flowchart of the Detailed Solution for the MHAM

3.2.2 Loss Function and Model Optimization

To guide the parameter learning of the entire model, a differentiable objective function that measures prediction bias needs to be defined. Since the task of this study is regression prediction, the traditional cross-entropy loss is no longer applicable.

This study constructs a composite loss function that includes mean squared error (MSE) and L_2 regularization to balance prediction accuracy and the model's generalization ability.

Main Prediction Loss: Mean Squared Error (MSE) To minimize the Euclidean distance between the model's predicted values and the true fault severity labels, mean squared error is used as the main loss term by (11).

$$L_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (11)$$

where N is the number of samples in the training batch, y_i is the true label of the i -th sample, and \hat{y}_i is the output value predicted by the model. This term forces the model output to approximate the real physical state as closely as possible.

To prevent the model from over-fitting in deep networks (especially multi-head attention layers with many parameters) and improve the robustness of the model under unknown conditions, this paper introduces an L_2 regularization term in the loss function by (12).

$$L_{Reg} = \sum_{w \in \Omega} w^2 \quad (12)$$

where Ω represents the set of all trainable weight parameters in the network. This term constrains the numerical range of the parameters, making the model smoother.

The final joint optimization objective function L_{Total} is defined as the weighted sum of the above two terms by (13).

$$L_{Total} = L_{MSE} + \frac{\lambda}{2} L_{Reg} \quad (13)$$

The optimal set of parameters θ^* is found to minimize the total loss by the goal of training the optimization model. The model is shown by (14).

$$\theta^* = \arg \min_{\theta} L_{Total} \quad (14)$$

This paper employs the Adam optimization algorithm for solving the problem. This algorithm dynamically adjusts the learning rate of each parameter using first- and second-order moment estimates of the gradient, effectively addressing the sparse gradient problem inherent in the MHAM and accelerating model convergence to the global optimum.

4 SIMULATION ANALYSIS

4.1 Experimental Setup and Data Description

To verify the effectiveness of the proposed method based on 1D-CNN-MHAM, the comparative simulation experiments are constructed based on the MATLAB (2023b) platform, with an AMD Ryzen 7 7735H processor and an NVIDIA RTX 4060 graphics card.

The data consisted of normalized multidimensional time-series signals. The dataset is randomly divided into training and test sets in an 8:2 ratio. The model is trained using the Adam optimizer with an initial learning rate of 0.001, a maximum number of iterations of 250, and a batch size of 128.

To demonstrate the algorithm's superiority, comparative experiments are conducted using 1D-CNN and the improved MHAM in this paper.

4.2 Evaluation Indicators

This study uses common quantitative indicators for four regression tasks to evaluate the predictive performance of the model, as shown in Table 2.

Table 2 Model Performance Evaluation Indicators

Evaluation indicators	Remark
RMSE	Reflects the model's ability to explain data variation.
MAE	Sample standard deviation of the deviation between predicted and actual values
MAPE	The actual average level reflecting the forecast error
R^2	Measure the percentage of prediction error relative to the true value.

4.3 Comparison and Analysis of Simulation Results

Intuitive Comparison of Fitting Performance

To visually demonstrate the predictive capabilities of the two models for the severity of braking system failures, line graphs comparing predicted and actual values on the test set, as well as scatter plots of linear regression fitting, are plotted for both models.

As shown in Figure 2, although the overall prediction trend of 1D-CNN follows the actual value, at extreme points of peaks and troughs, the predicted value often fails to reach the actual peak value, exhibiting a significant peak-shaving phenomenon. Its RMSE is 3.9662, indicating that under conditions of drastic signal fluctuations, it is difficult to accurately regress the peak value using only local convolutional features.

After introducing the attention mechanism, the fit between the predicted and actual curves is significantly enhanced, especially at local abrupt change points. The attention mechanism effectively compensates for the limitations of the local receptive field of the convolutional network through global weighting. Its RMSE further decreases to 3.9172, and the error fluctuation amplitude is significantly narrowed. Particularly in regions of signal abrupt change, the attention mechanism significantly corrects the bias of local convolutional features through weighting global contextual information, indicating that the model has higher robustness under complex nonlinear conditions.

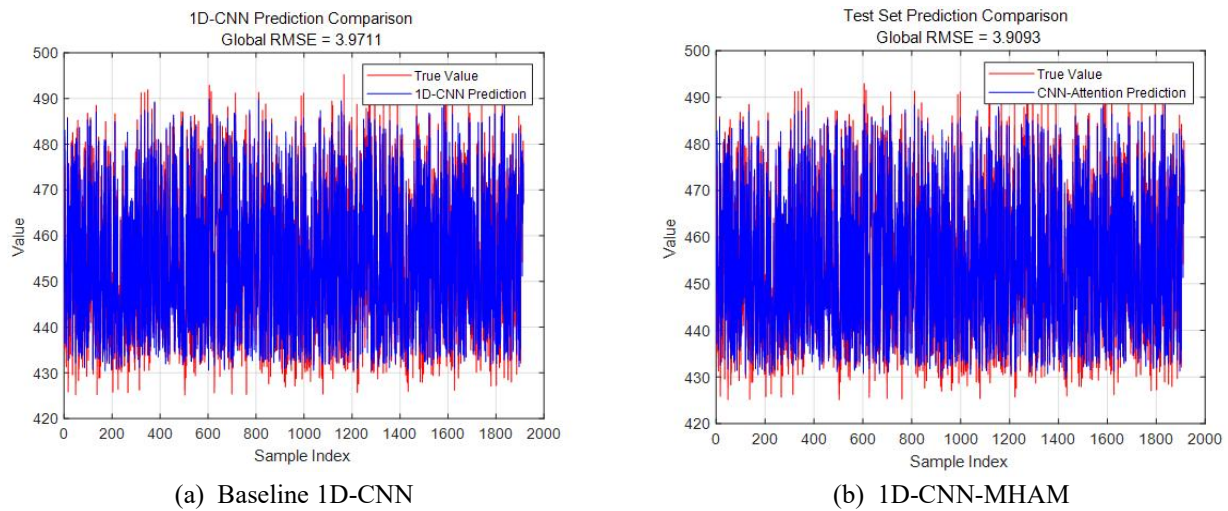


Figure 2 Comparison of RMSE between 1D-CNN and 1D-CNN-MHAM

This conclusion is further validated by correlation analysis, and regression correlation analysis (Figure 3) further confirmed it. Although the baseline model achieved a high degree of fit, significant outliers revealed its insufficient generalization to extreme samples. In contrast, the scatter plots of the improved model converged more closely to the ideal $y=x$ line, effectively suppressing outliers. This indicates that the model successfully broke through the performance bottleneck under high precision and uncovered deep and subtle correlations in the data.

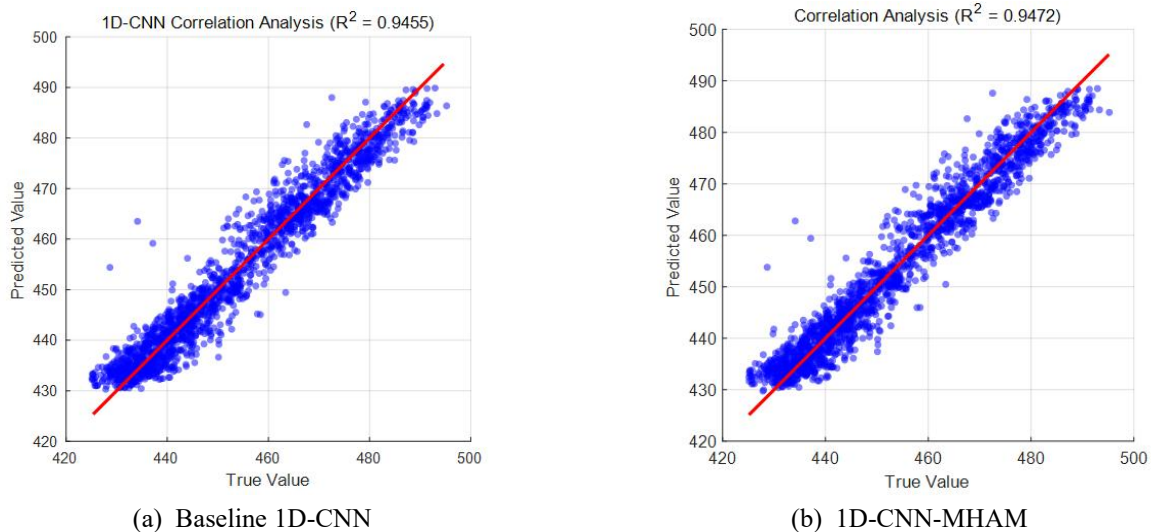


Figure 3 Comparison of Correlation between 1D-CNN and 1D-CNN-MHAM

4.4 Quantitative Analysis of Error Indicators

To more accurately quantify the performance differences between the models, this study statistically analyzed the key error metrics of the two algorithms on the test set, and the analysis data are shown in Table 3.

Table 3 Model Performance Comparison

	MAPE	MAE	MSE	RMSE	R^2
1D-CNN	0.68%	3.0846	15.7309	3.9711	0.9455
1D-CNN-MHAM	0.66%	3.0056	15.2829	3.9093	0.9472

4.5 Chapter Summary

This chapter uses comparative simulation experiments to systematically verify the performance of the fault diagnosis model based on 1D-CNN and the MHAM. Experimental results show that although the baseline 1D-CNN model already possesses good predictive capabilities, the introduction of the attention mechanism demonstrates advantages in suppressing signal noise, capturing long-range dependencies, and reducing extreme prediction errors. The 1.56% decrease in RMSE and the improved scatter plot convergence fully demonstrate the effectiveness of the proposed hybrid architecture and validate its engineering application value in improving the accuracy of fault diagnosis in active braking systems.

5 CONCLUSION

This paper addresses the challenge of quantifying fault characteristics in active braking systems under complex operating conditions by proposing a fault diagnosis method that integrates the 1D-CNN with the MHAM. This method efficiently integrates the local time-frequency features and global temporal dependencies of multi-source signals through a serial architecture, constructing an end-to-end quantitative assessment model of fault severity, providing accurate decision support for condition-based maintenance of the system.

Simulation results demonstrate that the proposed architecture exhibits significant advantages in regression tasks. Compared to traditional CNN networks, this architecture maintains high fitting accuracy while reducing RMSE and improving prediction accuracy. It confirms the effectiveness of the attention mechanism in suppressing noise interference and eliminating extreme prediction biases, significantly enhancing the model's robustness and generalization ability in nonlinear, strongly coupled environments.

In future research, this architecture provides a reference for lightweight deployment in edge devices. Furthermore, it offers concrete ideas for researching few-shot learning or transfer learning strategies to address the scarcity of data for certain extreme faults in real-world engineering scenarios, further improving the adaptability and reliability of diagnostic models in open and dynamic environments.

COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

REFERENCES

- [1] Lei Y, Yang B, Jiang X, et al. Applications of machine learning to machine fault diagnosis: A review and comparison. *Mechanical Systems and Signal Processing*, 2020, 138: 106587.
- [2] Zhao R, Yan R, Chen Z, et al. Deep learning and its applications to machine health monitoring. *Mechanical Systems and Signal Processing*, 2019, 115: 213-237.
- [3] Li Y, Wang X, Zhang Z. A review on deep learning in fault analysis of complex systems. In: *International Conference on Power and Energy Systems*, 2023.
- [4] Kiranyaz S, Avci O, Abdeljaber O, et al. 1D convolutional neural networks and applications: A survey. *Mechanical Systems and Signal Processing*, 2021, 151: 107398.
- [5] Eldele E, Ragab M, Chen Z, et al. TSLANet: Rethinking Transformers for Time Series Representation Learning. *arXiv preprint arXiv:2404.08472*, 2024.
- [6] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in Neural Information Processing Systems (NIPS)*, 2017: 5998-6008.
- [7] Li X, Zhang M, Guo H. Dual-Path Fault Diagnosis of Small Sample for Mechanical Systems Based on Multiple Attention Mechanisms. *IEEE Access*, 2024, 12: 114538-114551.
- [8] Jo T, Park I, Lee J, et al. A Fault Diagnosis and Fault-Tolerant Anti-Lock Brake System Control for Actuator Stuck Failures in Braking System in Autonomous Vehicles. *IEEE Transactions on Transportation Electrification*, 2025, 11(1): 188.
- [9] Li Y, Cheng J, Zhang W. Transformer network enhanced by dual convolutional neural network and cross-attention for wheelset bearing fault diagnosis. *Frontiers in Physics*, 2025, 13: 1546620.
- [10] Sifat M S I, Kabir M A, Islam M M M, et al. GAN-Based Data Augmentation for Fault Diagnosis and Prognosis of Rolling Bearings: A Literature Review. *IEEE Access*, 2025, 12: 1-21.
- [11] Kiranyaz S, Avci O, Abdeljaber O, et al. 1D convolutional neural networks and applications: A survey. *Mechanical Systems and Signal Processing*, 2021, 151: 107398.