

## Research Article

<https://doi.org/10.1631/jzus.A2400549>

# Self-attention and convolutional feature fusion for real-time intelligent fault detection of high-speed railway pantographs

Xufeng LI<sup>1</sup>, Jien MA<sup>1</sup>✉, Ping TAN<sup>2</sup>✉, Lanfen LIN<sup>3</sup>, Lin QIU<sup>1</sup>, Youtong FANG<sup>1</sup>

<sup>1</sup>College of Electrical Engineering, Zhejiang University, Hangzhou 310027, China

<sup>2</sup>School of Automation and Electricity Engineering, Zhejiang University of Science and Technology, Hangzhou 310023, China

<sup>3</sup>College of Computer Science and Technology, Zhejiang University, Hangzhou 310027, China

**Abstract:** Currently, most trains are equipped with dedicated cameras for capturing pantograph videos. Pantographs are core to the high-speed-railway pantograph-catenary system and their failure directly affects the normal operation of high-speed trains. However, given the complex and variable real-world operational conditions of high-speed railways, there is no real-time and robust pantograph fault-detection method capable of handling large volumes of surveillance video. Hence, it is of paramount importance to maintain real-time monitoring and analysis of pantographs. Our study presents a real-time intelligent detection technology for identifying faults in high-speed railway pantographs, utilizing a fusion of self-attention and convolution features. We delved into lightweight multi-scale feature-extraction and fault-detection models based on deep learning to detect pantograph anomalies. Compared with traditional methods, this approach achieves high recall and accuracy in pantograph recognition, accurately pinpointing issues like discharge sparks, pantograph horns, and carbon pantograph-slide malfunctions. After experimentation and validation with actual Electric Multiple Units train-surveillance videos, our algorithmic model demonstrates real-time, high-accuracy performance even under complex operational conditions.

**Key words:** High-speed railway pantograph; self-attention; CNN; real-time; feature fusion; fault detection

## 1 Introduction

China's high-speed railway has gradually become a network. Since the introduction of the 'Mid-Long Term Railway Network Plan' in 2016, 70% of the total operational scale of the 'eight vertical and eight horizontal' high-speed railway network has been completed (Tan et al., 2016). By 2023, the total operating mileage of China's high-speed railway had exceeded 45,000 kilometers. The amount of new high-speed railway network is tremendous, which creates new requirements and challenges for the daily maintenance of railway safety operations (Tan et al., 2020; Tan et al., 2021; Tan et al., 2022).

### Fig. 1 The pantograph-catenary system and overall model

The kinetic energy of a high-speed train is provided by the pantograph-catenary system, as shown in Fig 1. The pantograph obtains electric energy through sliding contact with the contact line and transmits it to the train to ensure normal operation (Zhou et al., 2011). Hence, high-quality current collection is a prerequisite for high-speed, safe, and reliable operation of trains. However, the contact position between the carbon slide and the contact line is the primary area of failure, and needs to be a focus of monitoring. In order to ensure safe operation of

✉ Jien MA, majien@zju.edu.cn

Ping TAN, 115011@zust.edu.cn

Xufeng LI, <https://orcid.org/0000-0002-0506-4834>

Jien MA, <https://orcid.org/0000-0003-0775-2793>

Ping TAN, <https://orcid.org/0000-0001-8656-3514>

Received Nov. 28, 2024; Revision accepted Jan. 3, 2025;  
Crosschecked

high-speed trains, the railway operation department must evaluate the state of key components of the arch network in a timely way. Due to the pantograph's structure and working characteristics, the non-contact monitoring method is safer. The most common method involves capturing pantograph surveillance video using a camera mounted on the train and transmitting it in real-time to the main system. At 350 km/h, any delay caused by pantograph failure can result in incalculable damage, with every second counting.

At present, the video monitoring system used for carriage pantographs only records and cannot perform real-time fault detection. In addition, the Auto Dropping Device (ADD) does not effectively detect pantograph faults or provide fall protection, a shortcoming which can lead to serious accidents. The intelligent analysis function based on surveillance video attempts to address this issue by detecting faults in time and thus enabling railway engineers to take relevant measures to avoid further deterioration of accidents. Thus, many researchers use image-processing technology to detect pantograph-catenary faults. Deng et al. (Deng et al., 2022) designed a pantograph slip-wear detection system using 3D structured light detection. Karaduman and Akin (Karaduman and Akin, 2022) proposed a new predictive maintenance method using the fuzzy classifier in railway systems. Wei et al. (Wei et al., 2020) proposed an innovative and intelligent method based on deep learning and image-processing technologies for online condition monitoring of the pantograph slide plate. Mo et al. (Mo et al., 2022) used the YOLOv4 (You Only Look Once) model, along with edge extraction and other traditional image-processing algorithms, to detect pantograph slide defects. Liu et al. (Liu et al., 2021) applied machine vision technology to measure the thickness of the slide, thereby reducing pantograph-catenary accidents caused by errors in manual measurement. Li et al. (Li et al., 2022) proposed an accurate real-time attitude-detection method for three-dimensional monocular pantographs which used the super ability of deep learning. Chen (Chen et al., 2022) proposed a depth-vision neural network detection method. Based on the YOLOv5 model, the depth pantograph-detection network (DPDN) was established to identify pantograph

regions in different complex scenes, and then the image visual feature extraction (IVFE) algorithm was used to detect the uneven distribution contact points between the pantograph and catenary in the pantograph region. Although various detection methods are available, significant challenges persist in complex real-world operational conditions, such as insufficient real-time capabilities and low fault-recognition rates.

As shown in Fig 1, the pantograph is divided into three regions: left, middle, and right. Although the middle region is the main working region, the failure of the left or right region will also cause the pantograph to lose balance, so detection of the pantograph cannot focus only on its contact-point region. Furthermore, the actual operating conditions of trains are complex, for example different bow types, different weather conditions, and trains passing over bridges and platforms, as shown in Fig. 2. Since the pantograph-monitoring scene is determined by actual operating conditions and the background is extremely complex, real-time processing is required in addition to efficient and robust image feature-extraction and classification methods. The integration of all these aspects is the main difficulty we hoped to solve in this study.

**Fig. 2 Pantograph detection using the proposed method in complex conditions**

Compared with traditional image-processing methods, deep convolutional neural networks have shown relatively strong performance. Since AlexNet (Krizhevsky et al., 2017) was proposed, deep convolutional neural networks have gradually become a hot research topic in the field of computer vision. In 2014, Simonyan K proposed VGGNet

(Visual Geometry Group Networks) (Simonyan and Zisserman, 2015), which has deep network structure, a small convolutional kernel, and a pooled sampling domain, enabling it to offer good transfer learning ability in feature extraction. In 2015, He et al. proposed ResNet (He et al., 2016) structure to solve the degradation problem caused by deep networks, and then proposed ResNetXt (Xie et al., 2017) and Resnet-D (He et al., 2019), which improved feature-extraction performance. In 2021, Ding et al. (Ding et al., 2021) proposed RepVGG, which only consisted of *Recall*,  $3 \times 3$  convolution, and ReLU (Rectified Linear Unit) activation functions, and further enhanced the feature-extraction performance of a VGG network through a simple branch-free structure. Although the depth and structure of convolutional networks are constantly improved and optimized, the characteristics of convolutional kernels limit the capacity for long-distance dependence. As Transformer (Vaswani et al., 2017) becomes a hot topic in natural language processing, more and more researchers are trying to transplant Transformer methods to image processing with long-term dependent features, in order to break the inherent limitations of convolutional networks. Vision Transformer (Dosovitskiy et al., 2021) realized application of Transformer in images for the first time. It uses image block coding as the network input and replaces convolutional operation with a self-attention module. It thus provides a new model in the field of computer vision, different from Convolutional Neural Networks (CNNs); but its performance needs to be improved. Considering the advantages of convolution, some researchers began to try composite models with both CNN and self-attention. For example, the CMT (Guo et al., 2022) model greatly improved feature-extraction performance. This indicates that the composite model has definite research value. To tackle fault detection in complex scenarios, a feature-extraction method with higher performance is required. Therefore, we chose to design a detection model here by combining Transformer and convolutional architectures.

To overcome the challenges identified above, we propose an inter-block feature fusion method. The contribution of the study is three-fold:

- The method integrates self-attention and convolutional features to improve

feature-extraction performance of convolutional networks and accurately identify pantographs in complex scenes.

- The lightweight multi-scale feature-extraction and fault-detection models are built to meet the requirements of real-time detection. Network parameters are reduced and model reasoning speed is improved.
- A complete and accurate fault-detection scheme for high-speed railway pantographs is established for the daily operation of trains.

## 2 Method

### 2.1 Overall model of pantograph detection

The overall model consists of two sub-models: a multi-scale feature-extraction model and a pantograph fault-detection model. In response to the limited number of samples of pantograph failures, we designed the method as follows: First, a feature-extraction network model is designed to learn the components of three regions using a large amount of normal pantograph image data. Then, a feature sample library of normal components is constructed. Finally, the presence of pantograph failures is detected by computing their confidence scores, as shown in Fig. 2. Compared to other types of faults, discharge spark faults have relatively simple shapes and colors, making them detectable through object recognition. As a result, these faults can also serve as subjects for learning. In this way, the problem of insufficient pantograph fault samples is avoided.

**Table 1 The backbone of the model**

Layer Number	Layer Name	Output Size	Parameter
E1	CBS	$320 \times 320$	$6 \times 6, 64, \text{stride } 2$
E2	CBS	$160 \times 160$	$3 \times 3, 128, \text{stride } 2$
E3	CBSBlock	$160 \times 160$	$[1 \times 1, 128] \times 3$
E4	CBS	$80 \times 80$	$3 \times 3, 256, \text{stride } 2$
E5	CBSBlock	$80 \times 80$	$[1 \times 1, 256] \times 6$
E6	CBS	$40 \times 40$	$3 \times 3, 512, \text{stride } 2$
E7	CBSBlock	$40 \times 40$	$[1 \times 1, 512] \times 9$
E8	CBS	$20 \times 20$	$3 \times 3, 1024, \text{stride } 2$
E9	SACNNBlock	$20 \times 20$	$[1 \times 1, 1024] \times 3$
E10	SACNN	$20 \times 20$	$[1 \times 1, 1024] \times 1$
E11	SPPF	$20 \times 20$	$5 \times 5, 1024, \text{pool}$

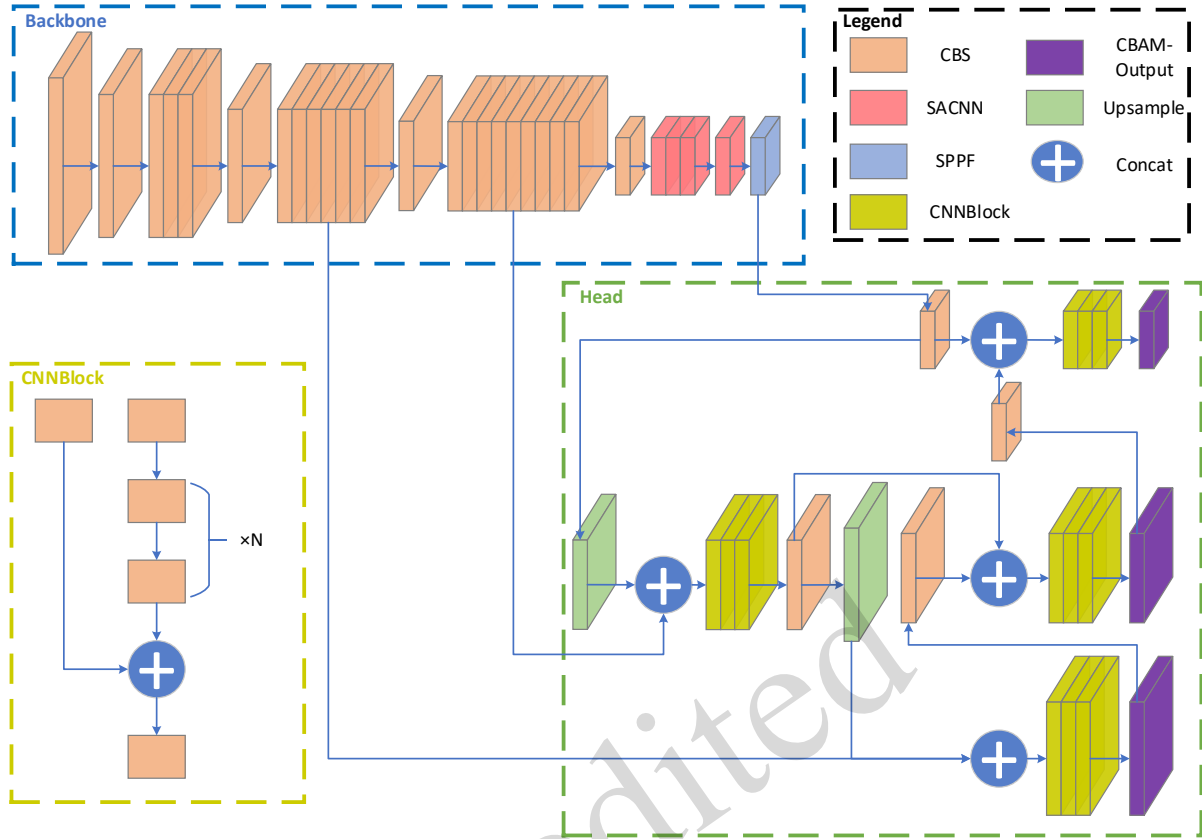


Fig. 3 Multi-scale feature-extraction model

Table 2 The head of the model

Layer Number	Layer Name	Output Size	Parameter
P1	CBS	20 × 20	1 × 1, 512
P2	Upsample	40 × 40	512
P3	Concat	40 × 40	1024, E7&P2
P4	CNNBlock	40 × 40	512
P5	CBS	40 × 40	1 × 1, 256
P6	Upsample	80 × 80	256
P7	Concat	80 × 80	512, E5&P6
P8	CNNBlock	80 × 80	256
P9	CBAM	80 × 80	256, Output1
P10	CBS	40 × 40	3 × 3, 256 stride 2
P11	Concat	40 × 40	512, P5&P10
P12	CNNBlock	40 × 40	512
P13	CBAM	40 × 40	512, Output2
P14	CBS	20 × 20	3 × 3, 512 stride 2
P15	Concat	20 × 20	1024, P1&P14
P16	CNNBlock	20 × 20	1024
P17	CBAM	20 × 20	1024, Output3

\*&: Concat.

The multi-scale feature-extraction model is shown in Fig. 3 and the specific parameters are shown in Tables 1 and 2, in which the CBS (Convolution and Batch Normalization and Sigmoid Linear Unit) module consists of three parts: convolution, batch standardization, and SiLU (Sigmoid Linear Unit) activation functions (Bochkovskiy et al., 2020). The function of BN is to make the network easier to converge. It also has a regularization function to prevent overfitting. Compared with ReLU, SiLU offers smoothness and non-monotonicity, which improves its effects. The calculation formula is given in Eq. (1).

$$f(x) = \frac{x}{1 + e^{-x}}. \tag{1}$$

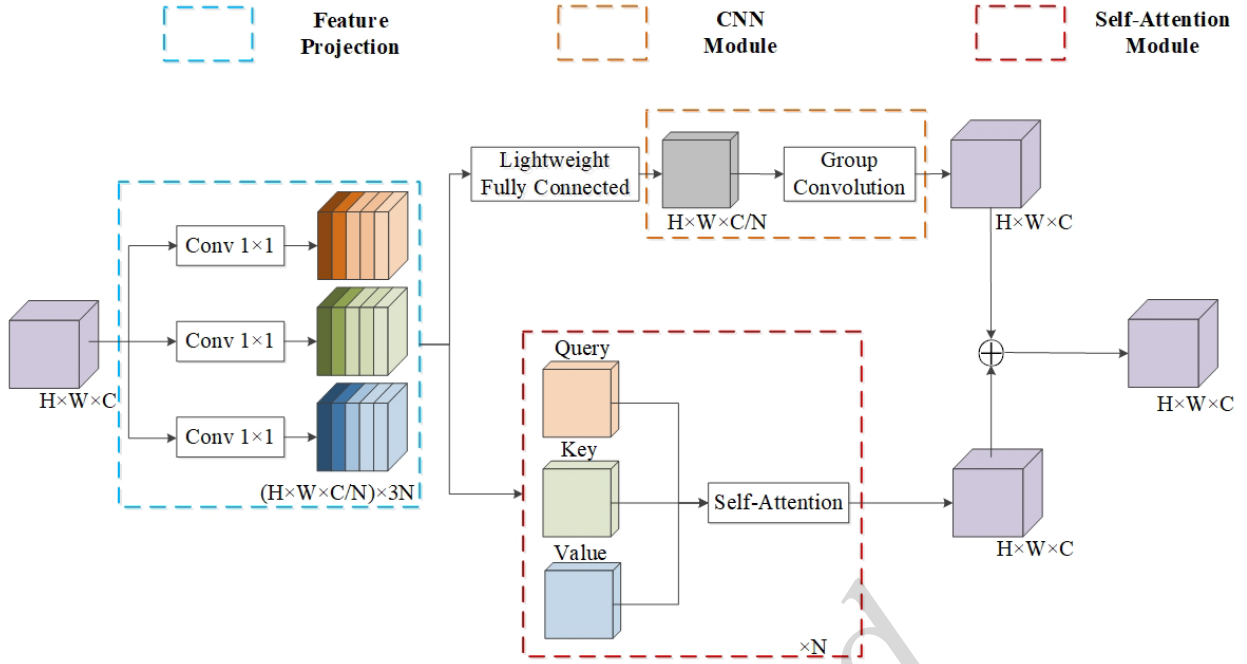


Fig. 4 A self-attention and convolutional neural network (SACNN)

To begin with, Table 1 shows a conventional feature-extraction network; this network is then connected with a FPN (Feature Pyramid Network) structure (Lin et al., 2017; Hao et al., 2022) and PAN (Pixel Aggregation Network) (Wang et al., 2019), as shown in Table 2. The FPN layer conveys strong semantic features from the top down (P1-P9), while the PAN conveys strong localization features from the bottom up (P10-P17). In this way, the parameters of different detection layers are aggregated from different backbone layers, further improving feature-extraction capability. The SACNNBlock (Self Attention and Convolutional Neural Network Block) module will be described in the next section. The CNNBlock module is similar to the C3 module in YOLO, as shown in Fig. 3. Similarly, the SPPF (Spatial Pyramid Pooling-Fast) (Bochkovskiy, et al., 2020) module at the bottom of the left network processes parallel inputs through multiple max-pools of varying sizes, followed by further fusion. This approach can address the multi-scale target problem to some extent. Finally, the network outputs three feature maps of different sizes for prediction, which can improve the recognition accuracy of targets of different sizes.

## 2.2 SACNNBlock

Similar to CSB superimposed on a CBSBlock module, the SACNNBlock module is overlaid by SACNN, which is formed by the integration of self-attention and CNN, as shown in Fig 4. It is an improved model based on ACmix (Pan et al., 2022). In the convolution path, it first uses a lightweight fully connected layer with a group number of  $N$  to aggregate the number of channels, and then restores the number of channels to  $C$  by a group convolution. By means of this, the amount of extra computation generated here is reduced from  $(3k^2 + k^4)C \times hw$  (Pan, et al., 2022) to  $(3 + k^2)C \times hw$  compared to ACmix.  $k$  is kernel size for convolution,  $C$  is input or output channel, and  $h$  and  $w$  are the length and width of the feature map, respectively.

First, we used three  $1 \times 1$  convolutional kernels for the convolutional operation of the input feature graphs, which are divided into  $N$  blocks to obtain  $3 \times N$  intermediate feature sets, and then we generated two calculation routes. The lower part of Fig. 4 shows the calculation route of self-attention, which generates three feature matrices, Query, Key, and Value, for the calculation of self-attention. The formula is given in Eq. (2) (Vaswani, et al., 2017; Dosovitskiy, et al., 2021; Hu et al., 2022).

$$\mathbf{Z} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_K}} + \mathbf{B}\right)\mathbf{V}, \quad (2)$$

where  $\mathbf{Z}$  is the calculation result of self-attention;  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$  are the three feature matrices of Query, Key, and Value respectively;  $d_K$  is the feature dimension of  $\mathbf{K}$ ; and  $\mathbf{B}$  is the relative position coding.

The upper part of Fig. 3 shows the convolutional calculation route, where  $3 \times N$  intermediate feature sets are spliced. After that, a lightweight fully connected layer is adopted to generate  $C/N$  feature maps, and a  $3 \times 3$  kernel is used for convolutional calculation. The obtained results are weighted and summed with the calculation result of self-attention. Finally, the calculation result of the fusion feature is output. It can be seen that the calculation process did not change the size of the feature map. The lightweight fully connected layer is composed of group convolution with the  $C/N$  group number, which is designed to reduce the number of parameters and the amount of computation.

## 2.2 Multi-scale feature output with CBAM

CBAM (Convolutional Block Attention Module) is a lightweight attention module that was created in 2018 and added to classical structures such as ResNet to improve the performance of convolutional networks (Woo et al., 2018; Wen et al., 2024). In this study, we added CBAM as a plug-in in front of three different output sizes for the network, as shown in Fig. 3.

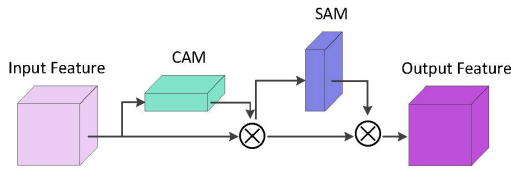


Fig. 5 CBAM

CBAM consists of two parts, namely a CAM (Channel Attention Module) and SAM (Spatial Attention Module) (Woo, et al., 2018; Ni et al., 2022; Yao et al., 2022), and its main structure is shown in Fig. 5.

### 2.2.1 CAM

First, the input feature images are obtained by max pooling and average pooling based on width and height, respectively, to obtain two  $1 \times 1 \times C$  feature images. Then, they are respectively sent into a two-layer MLP (Multilayer Perceptron) (Tan et al., 2024). The number of neurons in the first layer is  $C/r$  ( $r$  is the reduction rate), the activation function is Relu, and the number of neurons in the second layer is  $C$ . This two-layer neural network is shared. Then, the features output by MLP are added and activated by sigmoid to generate the final channel-attention feature. Finally, the channel attention feature is multiplied with the input feature to generate the input features required by SAM.

### 2.2.2 SAM

The output features of CAM are used for the input features of this module. First, two  $H \times W \times 1$  features are obtained by max pooling and average pooling based on channel, and then they are concatenated. After a  $7 \times 7$  convolutional operation, the dimension is reduced to 1 channel, namely  $H \times W \times 1$ . Next, a spatial attention feature is generated by sigmoid. Finally, this feature is multiplied with the input feature of this module to get the final output feature.

## 2.3 Loss function

The loss function in this study includes classification loss ( $\text{loss}_{\text{cls}}$ ), bounding-box regression loss ( $\text{loss}_{\text{box}}$ ), and objectness loss ( $\text{loss}_{\text{obj}}$ ), in which the  $\text{loss}_{\text{cls}}$  and  $\text{loss}_{\text{obj}}$  are calculated by the binary cross entropy loss function; the calculation formulas are given in Eqs. (3) and (4) (Bochkovski, et al., 2020).

$$\text{loss}_{\text{cls}} = \frac{1}{n} \sum_x [y^* \times \ln y + (1 - y^*) \times \ln(1 - y)], \quad (3)$$

$$\text{loss}_{\text{obj}} = \frac{1}{n} \sum_x [z^* \times \ln z + (1 - z^*) \times \ln(1 - z)], \quad (4)$$

where  $n$  is the total number of samples,  $x$  is samples,  $y^*$  and  $z^*$  are label values,  $y$  and  $z$  are predictive value, and  $\text{loss}_{\text{cls}}$  and  $\text{loss}_{\text{obj}}$  are classification loss and objectness loss, respectively.

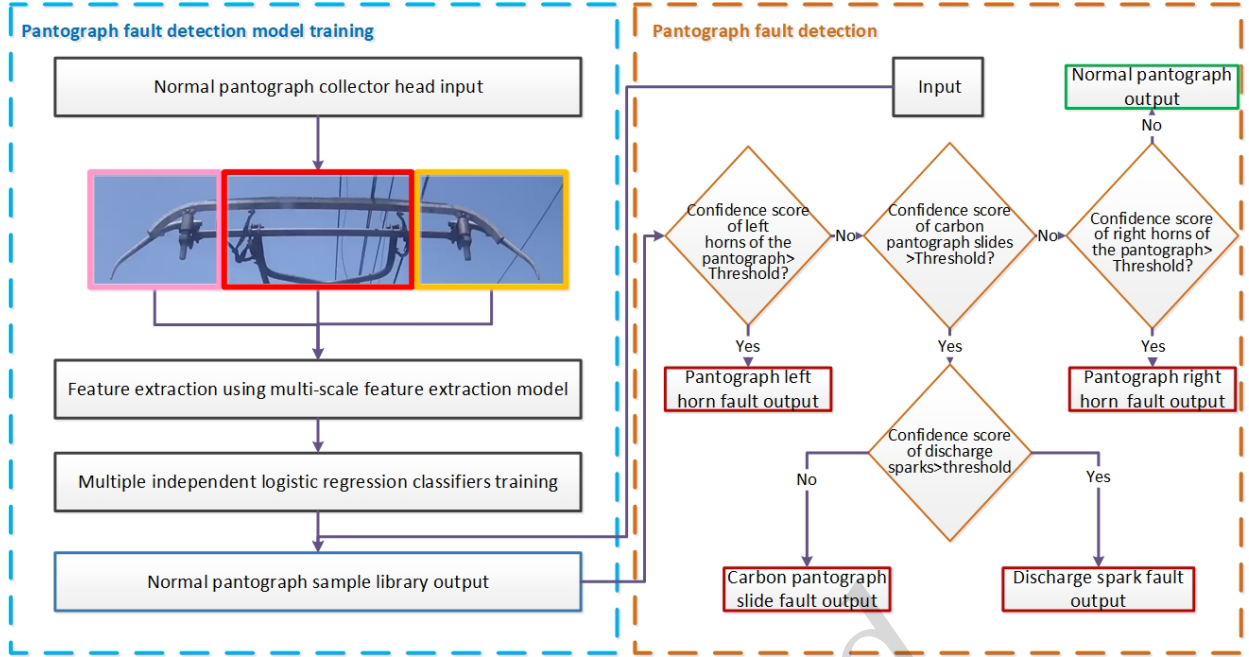


Fig. 6 Pantograph fault-detection model

CIoU loss (Complete Intersection over Union loss) is used for  $loss_{box}$ . Eqs. (5) - (8) give the calculation formulas (Zheng et al., 2020; Yao, et al., 2022).

$$loss_{box} = 1 - IoU + \frac{\rho^2(A, B)}{c^2} + \alpha v, \quad (5)$$

$$IoU = \frac{|A \cap B|}{|A \cup B|}, \quad (6)$$

$$v = \frac{4}{\pi^2} \left( \arctan \frac{w^B}{h^B} - \arctan \frac{w^A}{h^A} \right)^2, \quad (7)$$

$$\alpha = \frac{v}{1 - IoU + v}, \quad (8)$$

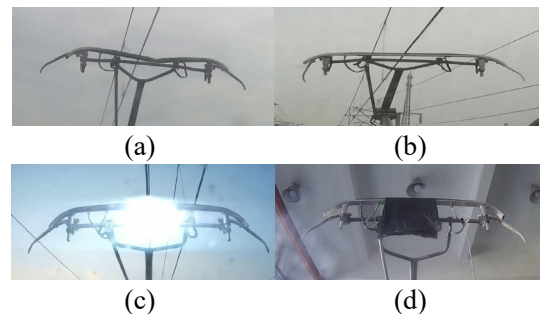
where  $loss_{box}$  is bounding-box regression loss,  $A$  is the predictive box,  $B$  is the target box,  $\rho$  is the distance between the center points of  $A$  and  $B$ ,  $c$  is the diagonal length of the minimum surrounding of  $A$  rectangle and  $B$  rectangle,  $v$  is the similarity of the aspect ratios of  $A$  and  $B$ ,  $\alpha$  is a positive trade-off parameter of  $v$ , and  $w$  and  $h$  are the width and height of the box respectively. The larger the IoU, the larger the overlap area between  $A$  and  $B$ , resulting in a larger  $\alpha$  and greater impact on  $v$ . In this case, it focuses on optimizing the aspect ratios of  $A$  and  $B$ . Otherwise, the smaller the IoU, the

smaller the overlap area between  $A$  and  $B$ , resulting in a smaller  $\alpha$  and smaller impact on the  $v$ . In this case, it focuses on optimizing the distance between  $A$  and  $B$ .

Therefore, the total loss is as follows:  $loss = \lambda_{cls} \times loss_{cls} + \lambda_{obj} \times loss_{obj} + \lambda_{box} \times loss_{box}$ , where  $\lambda$  is the gain coefficient of each loss function and the default value is 1.

#### 2.4 Pantograph fault-detection model

The pantograph image is divided into four parts: the pantograph collector head, the left horns of the pantograph, the carbon pantograph slide, and the right horns of the pantograph. There is also an extra part, the discharge spark. Only abnormal discharge sparks with large area are detected here. The detection process is shown in Fig. 6 and the pantograph fault types are shown in Fig. 7.



**Fig. 7 (a) Carbon pantograph slide deformation; (b) Horns of the pantograph fault; (c) Discharge spark; (d) Foreign body attachment**

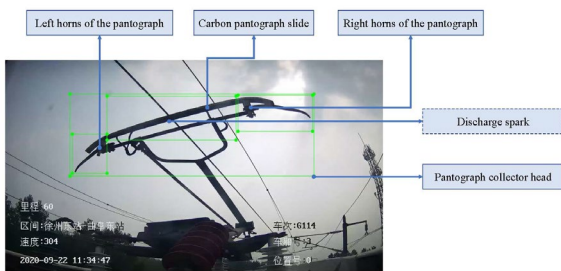
The blue dashed box in Fig. 6 illustrates the training process of the pantograph fault-detection model. It begins with the input of normal pantograph collector heads, extracted by the target-recognition network, which is subsequently divided into three regions and fed into a multi-scale feature-extraction model. Next, multiple independent logistic regression classifiers are utilized for training, resulting in the establishment of a feature sample library containing normal structures within these three regions.

The orange dashed box in Fig. 6 depicts the detection process. Here, the pantograph image under examination is inputted into the model. Its features are then compared with those in the standard library, generating a confidence score. This score is determined by multiplying the objectness with classification. Finally, these results are sequentially compared with the thresholds set within each region to ascertain the presence of a fault. If the score exceeds the threshold value, the object is considered normal; otherwise, it is deemed abnormal. The methodology for selecting these thresholds will be elucidated further in the subsequent chapter.

### 3 Experiments

#### 3.1 The pantograph dataset

There are a total of 2954 images in the training set, and each image sample is marked according to the standards in Fig. 8. The test set consisted of 600 pantograph images, including 500 normal pantograph images and 100 pantograph fault images.



**Fig. 8 Annotation of pantograph training images**

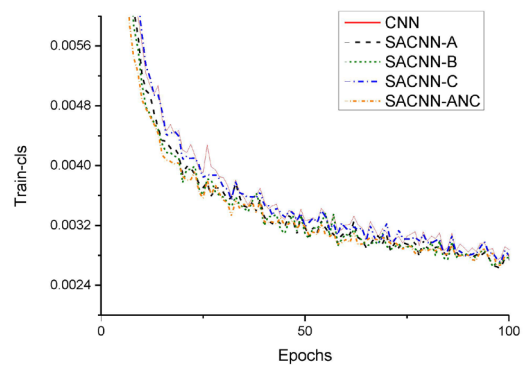
Five categories of images are marked: pantograph collector head, left horns of the pantograph, carbon pantograph slide, right horns of the pantograph, and discharge spark. We used GTX1050Ti for training and testing. The images size is set to  $640 \times 640$  and the batch size is set to 2; the number of training epochs is set to 200.

#### 3.2 Metrics

In this study, accurate fault detection is the most fundamental function, so the evaluation formula is as follows:  $\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$ ,  $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$ ,  $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$ , where TP means prediction is true and actual is true, TN means prediction is false and actual is false, FP means prediction is true and actual is false, and FN means prediction is false and actual is true. In addition, in order to better determine the confidence threshold, the measurement standard is introduced here, and the calculation formula is as follows:  $\text{F1} = 2 \times \text{Precision} \times \text{Recall} / (\text{Precision} + \text{Recall})$ .

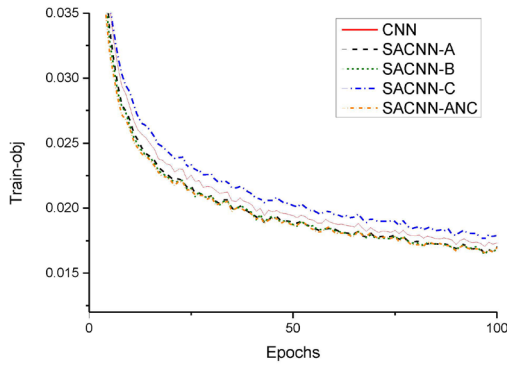
#### 3.3 Analysis of the training process

For each of the three branches of the loss function, the analysis of the training process is shown in Fig. 9.

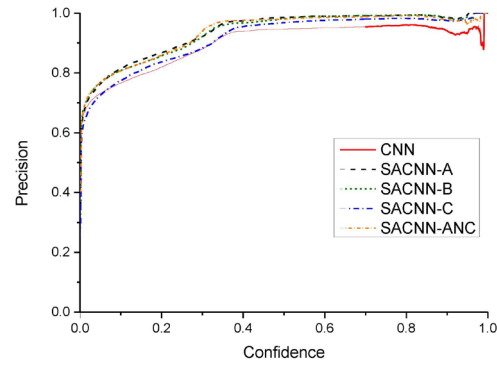


(a)

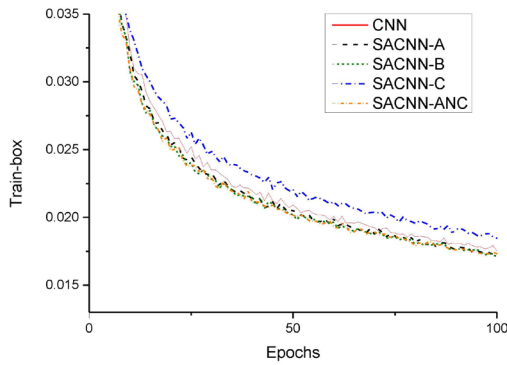




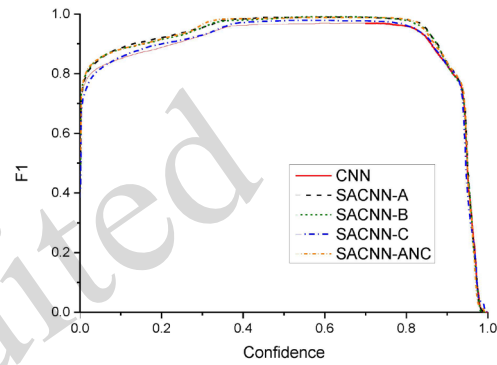
(b)



(b)



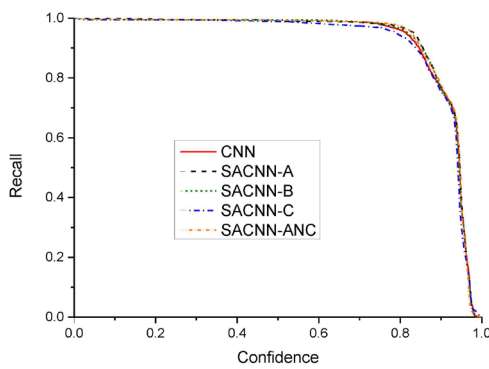
(c)



(c)

**Fig. 9** (a) Training loss of classification; (b) Training loss of objectness; (c) Training loss of bounding-box regression. CNN, multi-scale feature extraction model.

**Fig.10** (a) Recall and Confidence correlation curve in the training set; (b) Precision and Confidence correlation curve in the training set; (c) F1 and Confidence correlation curve in the training set



(a)

The CNN used is the multi-scale feature-extraction model, whereas all modules in E1-10 of Table 1's are convolutional modules. SACNN-A indicates that only the E10 module consists of the SACNN module. SACNN-ANC indicates that the CBAM is not used based on SACNN-A. SACNN-B indicates that only the layer E9 consists of SACNNBlock module and layer E10 consists of SACNN module. SACNN-C indicates that only the layer E7 and E9 consist of SACNNBlock module and layer E10 consists of SACNN module. As can be seen from the diagram of the training process (Fig. \*), in  $loss_{cls}$ , there is little difference between CNN and SACNN-C, while the loss curves of SACNN-A and SACNN-B are all lower than the former, indicating lower loss. Moreover, on  $loss_{obj}$  and  $loss_{box}$ , when a SACNN module is added to E9 and E10 layers, the training effect is better than that of CNN, but when

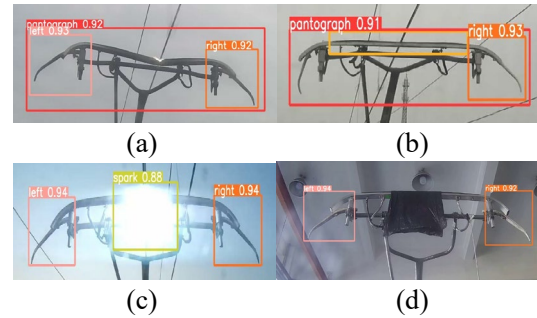
SACNN is added to the E7 layer, the training effect is worse. It can be concluded that when self-attention acts on the shallow layer of the network, the effect is worse than that of CNN, but when self-attention acts on the deep layer of the network, it can better improve the network performance of CNN.

### 3.4 Influence of the number of SACNN modules and the confidence-score threshold

The higher the confidence, the higher the inference accuracy of the object. As can be seen from Fig. 10, when the confidence level is 0.8, Recall and Precision of all classification averages in the training set, as well as the confidence, are relatively high. Moreover, from the F1 curve of integrated Recall and Precision, it is evident that when the confidence value exceeded 0.8, model performance began to decline. Therefore, the confidence of object detection is set at 0.8 as the threshold value in the fault-detection model.

Table 3 gives the test results of the pantograph fault detection. When the SACNN module is added, the accuracy of the models is better than that of CNN, but the reasoning speed of the models decreased as more modules are added, and the memory also increased. In the test of these four models, SACNN-A successfully detected all pantograph failures, and its overall performance is significantly better than that of the other models; its speed also met the requirements of real-time detection. This proved that adding a CBAM module before network output is helpful in improving the accuracy of object recognition.

### 3.5 Pantograph-detection model test results



**Fig. 11** Pantograph fault-detection results. (a) Carbon pantograph slide deformation; (b) Horns of the pantograph fault; (c) Discharge spark; (d) Foreign body attachment

Fig. 11 shows the detection results for several types of pantograph faults. It is clear that when there is a fault in a certain region of the pantograph, the object in the region could not be identified normally, so we determined that there is a fault in the region. Fig. 11 (c) demonstrates that this method not only detected faults in the carbon slider but also identified the specific category of major spark faults. Fig. 11 (d) shows foreign bodies attached to the carbon pantograph slide. Due to the wide variety of foreign bodies, direct identification is difficult to achieve, so the method proposed here instead aims to identify faults of the carbon pantograph slide in this region. Adopting the approach outlined by Chen, et al. (2022) would be ineffective for identifying the two types of faults depicted in (a) and (b). In Fig. 11 (a), despite deformation of the carbon slider, its contact point remains within the working area. Similarly, in Fig. 11 (b), the fracture at the left bow angle does not affect the contact point.

**Table 3** Pantograph detection model test results with different numbers of SACNN modules

Model	Recall (%)	Precision (%)	Accuracy (%)	FPS	Memory(G)
CNN+PFDM*	93.0	63.7	90.0	66.67	<b>0.562</b>
SACNN-A+PFDM	<b>100.0</b>	85.2	<b>97.1</b>	62.50	0.686
SACNN-B+PFDM	89.0	77.3	68.8	32.25	0.797
SACNN-C+PFDM	93.0	68.8	91.8	16.13	1.660
SACNN-ANC+PFDM	93.0	<b>85.9</b>	96.3	<b>66.96</b>	0.682

\*PFDM: Pantograph Fault-Detection Model.

**Table 4** Comparison with other methods of pantograph fault detection

Model	Recall (%)	Precision (%)	Accuracy (%)	FPS	Memory(G)
SSD+PFDM*	82.0	65.2	89.7	5.88	2.581
Faster-RCNN+PFDM	89.0	<b>99.5</b>	<b>98.1</b>	5.55	2.581
YOLOv3+PFDM	82.0	90.1	95.5	14.08	2.330

YOLOv5+PFDM	88.0	88.0	96.0	40.00	1.740
YOLOv7+PFDM	<b>100.0</b>	34.6	68.5	8.11	1.980
ACmix+PFDM	94.0	81.7	95.5	56.38	0.686
DPDN+IVFE	43.0	11.3	34.5	22.22	1.788
Proposed method	<b>100.0</b>	85.2	97.1	<b>62.50</b>	<b>0.686</b>

We then compared the model horizontally with other methods, and the results are shown in Table 4. In this test, the SACNN module of SACNN-A is replaced by an ACmix module to build an ACmix detection model. Based on the experimental results from Chen, et al. (2022) in Table 4, it is evident that the comprehensive accuracy of a fault-detection method decreases significantly when the dataset is extended to encompass complex scenarios involving all-weather conditions and various shooting perspectives. Under the same conditions, the PFDM (Pantograph Fault-Detection Model) exhibited higher accuracy. Furthermore, compared to other object feature detection methods, it can be clearly seen that the performance of the proposed model is superior.

When the graphics card is configured as GTX1050Ti, the FPS (Frames Per Second) of the PFDM is 62.50, while the FPS of images captured by the surveillance camera on the high-speed railway is 25, which meets the requirements of real-time performance. Although the Faster-RCNN method can achieve higher accuracy, its FPS is lower than the real-time requirement of 25. At the same time, the proposed method needs less memory and is lighter in weight, so it is convenient to use in high-speed trains for real-time detection of pantographs.

## 4 Conclusions

In order to solve the problems involved in obtaining real-time intelligence on high-speed railway pantograph faults, we propose a lightweight deep learning model based on the fusion of self-attention and convolutional features to achieve real-time and high-accuracy recognition of key components of pantographs. Our experiments show that the fusion module can effectively improve the performance of the original convolutional network. High recall and accuracy are achieved in training and testing sets, and the algorithm model has good performance. We designed a pantograph fault-detection model which achieves fast and intelligent fault detection and ac-

curately identified all the faults in the test set. The detection model proposed here provides a set of real-time and accurate tools for intelligent fault detection in pantographs of high-speed trains, and provides an efficient auxiliary means for engineers on board to assess pantograph faults.

## Acknowledgments

This work was supported by National Key R&D Program of China under Grant (2022YFB4301102).

## Author contributions

Xufeng LI designed the research. Lin QIU and Youtong FANG processed the corresponding data. Xufeng LI wrote the first draft of the manuscript. Ping TAN and Lanfen LIN helped to organize the manuscript. Jien MA revised and edited the final version.

## Conflict of interest

Xufeng LI, Jien MA, Ping TAN, Lanfen LIN, Lin QIU and Youtong FANG declare that they have no conflict of interest.

## References

- Bochkovski A, Wang C-Y, Liao H-YM, 2020. Yolov4: Optimal speed and accuracy of object detection. *CoRR*, abs/2004.10934
- Chen RC, Lin YZ, Jin T, 2022. High-speed railway pantograph-catenary anomaly detection method based on depth vision neural network. *Ieee Transactions on Instrumentation and Measurement*, 71 <https://doi.org/10.1109/tim.2022.3188042>
- Deng Q, Chen U, Yipin Z, 2022. Research on wear detection of pantograph slide plate based on high speed 3d structured light detection. *Engineering and Technological Research*, 7(12):15-19. <https://doi.org/10.19537/j.cnki.2096-2789.2022.12.005>
- Ding XH, Zhang XY, Ma NN, et al., 2021. Repvgg: Making vgg-style convnets great again. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Electr Network. p.13728-13737. <https://doi.org/10.1109/cvpr46437.2021.01352>
- Dosovitskiy A, Beyer L, Kolesnikov A, et al., 2021. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*,
- Guo JY, Han K, Wu H, et al., 2022. Cmt: Convolutional neural

- networks meet vision transformers. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA. p.12165-12175. <https://doi.org/10.1109/cvpr52688.2022.01186>
- Hao K, Chen GK, Zhao L, et al., 2022. An insulator defect detection model in aerial images based on multiscale feature pyramid network. *Ieee Transactions on Instrumentation and Measurement*, 71 <https://doi.org/10.1109/tim.2022.3200861>
- He KM, Zhang XY, Ren SQ, et al., 2016. Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Seattle, WA. p.770-778. <https://doi.org/10.1109/cvpr.2016.90>
- He T, Zhang Z, Zhang H, et al., 2019. Bag of tricks for image classification with convolutional neural networks. 32nd IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA. p.558-567. <https://doi.org/10.1109/cvpr.2019.00065>
- Hu TY, Ma HM, Liu H, et al., 2022. Self-attention-based machine theory of mind for electric vehicle charging demand forecast. *Ieee Transactions on Industrial Informatics*, 18(11):8191-8202. <https://doi.org/10.1109/tii.2022.3180399>
- Karaduman G, Akin E, 2022. A new approach based on predictive maintenance using the fuzzy classifier in pantograph-catenary systems. *Ieee Transactions on Intelligent Transportation Systems*, 23(5):4236-4246. <https://doi.org/10.1109/tits.2020.3042997>
- Krizhevsky A, Sutskever I, Hinton GE, 2017. Imagenet classification with deep convolutional neural networks. *Communications of the Acm*, 60(6):84-90. <https://doi.org/10.1145/3065386>
- Li D, Pan X, Fu ZZ, et al., 2022. Real-time accurate deep learning-based edge detection for 3-d pantograph pose status inspection. *Ieee Transactions on Instrumentation and Measurement*, 71 <https://doi.org/10.1109/tim.2021.3137558>
- Lin TY, Dollár P, Girshick R, et al., 2017. Feature pyramid networks for object detection. 30th IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI. p.936-944. <https://doi.org/10.1109/cvpr.2017.106>
- Liu W, Zhao J, Wang S, 2021. Pantograph slide thickness detection method research based on machine vision. *Electronic Measurement Technology*, 44(24):128-133. <https://doi.org/10.19651/j.cnki.emt.2107928>
- Mo X, Wang K, Pan C, et al., 2022. Intelligent defect-detection technology of pantograph pan based on the image from railway 5c device. *China Railway*, (02):148-155. <https://doi.org/10.19549/j.issn.1001-683x.2021.07.14.002>
- Ni XF, Ma ZJ, Liu JW, et al., 2022. Attention network for rail surface defect detection via consistency of intersection-over-union(iou)-guided center-point estimation. *Ieee Transactions on Industrial Informatics*, 18(3):1694-1705. <https://doi.org/10.1109/tii.2021.3085848>
- Pan XR, Ge CJ, Lu R, et al., 2022. On the integration of self-attention and convolution. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA. p.805-815. <https://doi.org/10.1109/cvpr52688.2022.00089>
- Simonyan K, Zisserman A, 2015. Very deep convolutional networks for large-scale image recognition. 3rd International Conference on Learning Representations, {ICLR} 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings,
- Tan B, Wang D, Shi J, et al., 2024. Temperature field prediction of steel-concrete composite decks using tvfemd-stacking ensemble algorithm. *Journal of Zhejiang University SCIENCE A*, 25(9):732-748.
- Tan P, Ma JE, Zhou J, et al., 2016. Sustainability development strategy of china's high speed rail. *Journal of Zhejiang University-Science A*, 17(12):923-932. <https://doi.org/10.1631/jzus.A1600747>
- Tan P, Li XF, Xu JM, et al., 2020. Catenary insulator defect detection based on contour features and gray similarity matching. *Journal of Zhejiang University-Science A*, 21(1):64-73. <https://doi.org/10.1631/jzus.A1900341>
- Tan P, Li XF, Wu ZG, et al., 2021. Multialgorithm fusion image processing for high speed railway dropper failure-defect detection. *Ieee Transactions on Systems Man Cybernetics-Systems*, 51(7):4466-4478. <https://doi.org/10.1109/tsmc.2019.2938684>
- Tan P, Li XF, Ding J, et al., 2022. Mask r-cnn and multifeature clustering model for catenary insulator recognition and defect detection. *Journal of Zhejiang University-Science A*, 23(9):745-756. <https://doi.org/10.1631/jzus.A2100494>
- Vaswani A, Shazeer N, Parmar N, et al., 2017. Attention is all you need. 31st Annual Conference on Neural Information Processing Systems (NIPS), Long Beach, CA.
- Wang WH, Xie EZ, Song XG, et al., 2019. Efficient and accurate arbitrary-shaped text detection with pixel aggregation network. IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, SOUTH KOREA. p.8439-8448. <https://doi.org/10.1109/iccv.2019.00853>
- Wei XK, Jiang SY, Li Y, et al., 2020. Defect detection of pantograph slide based on deep learning and image processing technology. *Ieee Transactions on Intelligent Transportation Systems*, 21(3):947-958. <https://doi.org/10.1109/tits.2019.2900385>
- Wen YT, Cheng JX, Ren YX, et al., 2024. Complex defects detection of 3-d-printed lattice structures: Accuracy and scale improvement in yolo v7. *Ieee Transactions on Instrumentation and Measurement*, 73 <https://doi.org/10.1109/tim.2024.3370765>
- Woo SH, Park J, Lee JY, et al., 2018. Cbam: Convolutional block attention module. 15th European Conference on Computer Vision (ECCV), Munich, GERMANY. p.3-19.

- [https://doi.org/10.1007/978-3-030-01234-2\\_1](https://doi.org/10.1007/978-3-030-01234-2_1)  
Xie SN, Girshick R, Dollár P, et al., 2017. Aggregated residual transformations for deep neural networks. 30th IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI. p.5987-5995. <https://doi.org/10.1109/cvpr.2017.634>
- Yao H, Liu YH, Li X, et al., 2022. A detection method for pavement cracks combining object detection and attention mechanism. *Ieee Transactions on Intelligent Transportation Systems*, 23(11):22179-22189. <https://doi.org/10.1109/tits.2022.3177210>
- Zheng ZH, Wang P, Liu W, et al., 2020. Distance-iou loss: Faster and better learning for bounding box regression. 34th AAAI Conference on Artificial Intelligence / 32nd Innovative Applications of Artificial Intelligence Conference / 10th AAAI Symposium on Educational Advances in Artificial Intelligence, New York, NY. p.12993-13000.
- Zhou N, Zhang WH, Li RP, 2011. Dynamic performance of a pantograph-catenary system with the consideration of the appearance characteristics of contact surfaces. *Journal of Zhejiang University-Science A*, 12(12):913-920. <https://doi.org/10.1631/jzus.A11GT015>

## 中文概要

**题目:** 基于自注意力和卷积特征融合的高铁受电弓实时智能故障检测

**作者:** 李旭峰<sup>1</sup>, 马吉恩<sup>1</sup>, 谭平<sup>2</sup>, 林兰芬<sup>3</sup>, 邱麟<sup>1</sup>, 方攸同<sup>1</sup>

**机构:** <sup>1</sup>浙江大学, 电气工程学院, 中国杭州, 310027;  
<sup>2</sup>浙江科技大学, 自动化与电气工程学院, 310023;  
<sup>3</sup>浙江大学, 计算机科学与技术学院, 中国杭州, 310027

**目的:** 由于高铁运行工况复杂, 对受电弓故障的实时检测监测技术存在较大难点。本文旨在基于受电弓监控视频, 研究实时智能故障检测方法, 及时发现故障, 保障列车安全运行。

**创新点:** 1. 将自注意力和卷积特征相结合, 提高了卷积网络的特征提取性能, 使其在复杂场景中准确识别受电弓; 2. 构建轻量级的多尺度特征提取和故障检测模型, 满足实时检测的要求。减少了网络参数, 提高了模型推理速度; 3. 针对列车的日常运行, 建立了一套完整、准确的高速铁路受电弓故障检测方案。

**方法:** 整个模型由两个子模型组成: 多尺度特征提取模型和受电弓故障检测模型。针对受电弓故障样本

数量少的问题, 该方法设计如下: 首先, 设计轻量化多尺度特征提取网络模型, 利用大量的正常受电弓视频图像数据学习各部件的特征。然后, 构建正常受电弓部件的特征样本库; 最后, 通过匹配正常样本库计算其置信度来检测受电弓故障。

**结论:** 1. 本文提出了一种基于自注意力特征与卷积特征融合的轻量级深度学习模型, 实现了对受电弓关键部件的实时、高精度识别。2. 实验表明, 融合模块可以有效地提高原卷积网络的性能。在训练集和测试集上均取得了较高的查全率和查准率, 算法模型具有良好的性能。3. 设计了受电弓故障检测模型, 实现了快速、智能的故障检测, 能够准确识别出测试集中的所有故障。

**关键词:** 高速铁路受电弓; 自注意力; 卷积神经网络; 实时; 特征融合; 故障检测