


RESEARCH

Open Access



MFVT: an anomaly traffic detection method merging feature fusion network and vision transformer architecture

Ming Li¹, Dezhi Han^{1*} , Dun Li¹, Han Liu¹ and Chin-Chen Chang²

*Correspondence:

dezhihan88@sina.com

¹ Department of Engineering,
Shanghai Maritime University,
Shanghai, China

Full list of author information
is available at the end of the
article

Abstract

Network intrusion detection, which takes the extraction and analysis of network traffic features as the main method, plays a vital role in network security protection. The current network traffic feature extraction and analysis for network intrusion detection mostly uses deep learning algorithms. Currently, deep learning requires a lot of training resources and has weak processing capabilities for imbalanced datasets. In this paper, a deep learning model (MFVT) based on feature fusion network and vision transformer architecture is proposed, which improves the processing ability of imbalanced datasets and reduces the sample data resources needed for training. Besides, to improve the traditional raw traffic features extraction methods, a new raw traffic features extraction method (CRP) is proposed, and the CRP uses PCA algorithm to reduce all the processed digital traffic features to the specified dimension. On the IDS 2017 dataset and the IDS 2012 dataset, the ablation experiments show that the performance of the proposed MFVT model is significantly better than other network intrusion detection models, and the detection accuracy can reach the state-of-the-art level. And, when MFVT model is combined with CRP algorithm, the detection accuracy is further improved to 99.99%.

Keywords: Network intrusion detection, Traffic features, Deep learning, Feature fusion network, Vision transformer, MFVT, CRP, Detection accuracy

1 Introduction

The rapid development of the mobile Internet not only brings great convenience to network users and society but also allows criminals to create a series of attacks in the network. These attacks have seriously threatened the normal operation of the network, not only caused a lot of economic losses but also brought hidden dangers to national security [1–3]. A group of behaviors that violate computer security policies such as confidentiality, integrity, and availability are defined as intrusion detection [4, 5]. As a security protection system used to monitor computer network, the intrusion detection system can detect suspicious behaviors and take corresponding measures to ensure the normal operation of the network and reduce economic losses, which has been in use since the 1980s [6, 7]. Recently, due to the rapid development of mobile Internet, attacks on Internet-connected devices are gradually increasing. Thus, many scholars have a strong

interest in the research of intrusion detection systems and good detection results have been achieved [8–10].

Besides, the detection of anomaly network traffic is an important task of network intrusion detection, which is essential to classify network traffics [11], which requires researchers to make accurate judgments on the collected network traffic data and detect network traffic with offensive behavior. To detect anomaly traffics more effectively, network traffic packets are usually divided into flows according to source IP, destination IP, source port, destination port, protocol, and timestamp [12]. The current anomaly traffic detection technology mainly includes traditional network anomaly traffic detection technology and network anomaly traffic detection method based on machine learning. In this paper, deep learning methods were used to classify network traffics. Deep learning methods have the characteristics of end-to-end and automatic extraction of network traffic data features, to avoid the cumbersome process of manual extraction of features, and deep learning methods have good adaptability, self-organization, and promotion ability. So, the use of deep learning can make the detection system have more stable performance and higher detection efficiency [13, 14].

However, deep learning technology needs a large amount of labeled data for training, and labeled data require experts with specific knowledge to spend a lot of time on labeling, which is time-consuming and laborious. Most of the datasets used in deep learning are imbalanced datasets. These problems cause a significant impact on the performance of deep learning models. Under-sampling and over-sampling are commonly used to solve data imbalance problems, but under-sampling will discard some data leading to the loss of some features, and over-sampling will add some data leading to changing the original data distribution, both of which have an impact on the experimental accuracy [15]. In this paper, the traffic features learned from a two-layer convolutional networks are fused, which can alleviate the impact of data imbalance on the accuracy of the experiment. Due to the outstanding performance of transformer architecture in the field of natural language processing (NLP) and the limitations of its application in computer vision, Dosovitskiy [16] improved the transformer architecture and proposed vision transformer architecture for image sequence converter realize image classification and achieved good results. Meanwhile, experiments proved that vision transformer required fewer training resources. Inspired by the vision transformer architecture, a deep learning model (MFVT) based on the feature fusion network and the vision transformer architecture was proposed in this paper for network anomaly traffic detection. MFVT model has strong ability to deal with imbalanced datasets and therefore effectively reduce the sample resources required for training. This paper also studies the influence of learning rate change and the number of training epochs on the experimental accuracy based on the MFVT model.

So far, there are many ways to process raw network traffic data, but there is no uniform standard. Since the data that a neural network can accept must be of the same dimension, the extracted network traffic data must be filtered to a specific dimension before it can be used as the input of the neural network model. Most of the traditional methods directly intercept the data of specific dimensions from the network traffic data. Although the effect is quite good, there is room for improvement. Therefore, PCA algorithm is used in this paper to reduce all the processed digital traffic features to a specified

dimension. The experimental accuracy obtained in the datasets IDS 2017 [17] and IDS 2012 [18] is significantly higher than the traditional methods.

In summary, the main contributions of this paper are as follows.

- A deep learning model (MFVT) based on feature fusion network and vision transformer architecture is proposed, which can effectively improve the detection accuracy while reducing the training resources. On the IDS 2017 dataset and the IDS 2012 dataset, MFVT model can achieve the best performance on all evaluation metrics.
- A new raw traffic data extraction algorithm (CRP) is proposed, which uses the PCA [19] algorithm to reduce the processed digital traffic features to a specified dimension. The ablation experiment results show that the detection accuracy has significantly improved to compare with traditional methods.
- Based on the MFVT model, the impact of training epochs and the variation of the learning rate on the detection performance of the model is further studied.

The rest of this paper is organized as follows. Section 2 introduces the related works to the model and method presented in this paper; Sect.3 details the deep learning model and the raw network traffic data processing algorithm, Sect.4 introduces ablation experiments and experimental results of MFVT model in detail, and finally, our work is summarized in Sect.5 .

2 Related work

This section mainly summarizes some documents related to the work of this paper, including intrusion detection and transformer architecture.

2.1 Intrusion detection

With the continuous development of artificial intelligence big data and cloud computing technology, intrusion detection technology is constantly updated using new technologies [20–23] In 1980, Anderson [24] proposed the concept of intrusion detection technology, which aims to timely identify abnormal behaviors in the network and reduce losses caused by abnormal behaviors. Over the past 40 years, many methods have been used in intrusion detection, all of which aim to sense attacks with good predictive accuracy and improve real-time prediction. These methods all attempting to extract a pattern from network traffics to distinguish attack traffics from regular traffics.

Specifically, Table 1 briefly summarizes the methods used in intrusion detection. Currently, the traditional machine learning methods applied to the field of intrusion detection are mainly supervised learning, such as support vector machine (SVM) [25–27], K-nearest neighbor (KNN) [28], and random forest (RF) [29, 30]. These methods mentioned above have a high false alarm rate and a low detection rate for attack traffics. It is a common problem in traditional machine learning methods to design a feature set that can accurately reflect traffic characteristics, and the quality of feature set directly affects the classification performance of the method. In recent years, although many researchers have been working on the problem of how to design feature sets [31, 32], how to design a set of suitable traffic feature sets is still an unresolved research topic.

Table 1 A brief summary of intrusion detection methods

Author	Method		DataSets	References
Yin C L	Machine learning	SVM	NSL-KDD	[25]
Reddy R.R			KDD99	[27]
Li W		KNN	Flooding Attack	[28]
Farnaaz N		RF	NSL-KDD Dataset	[29]
Zhang J			KDD99	[30]
Yan Q	Deep learning	CNN	KDDCUP'99	[34]
Zhang Y		CNN_LSTM	CICIDS2017	[13]
Lin P		Attention+LSTM	CSE-CIC-IDS2018	[35]
Zhang Y		PCCN	CICIDS2017	[36]
Zhong Y		HELAD	KDDCUP99 +CICIDS2017	[37]

Moreover, deep learning [33] has good self-adaptability, self-organization, and generalization capabilities. Therefore, it can be a good solution to the problem that traditional machine learning needs to manually design a group of feature sets. The use of deep learning can enable detection systems with higher detection efficiency and therefore has been widely studied by scholars in recent years. Yan [34] constructed an intrusion detection system based on convolutional neural network (CNN) and applied generative adversarial network to synthesize attack traces, and experimental results verified the effectiveness of the system. Zhang [13] proposed a deep hierarchical network-based intrusion detection model that combines CNN and long short-term memory (CNN_LSTM) network, and the CNN_LSTM model achieved good performance on the IDS2017 dataset. Lin [35] constructed a dynamic network anomaly detection system, which uses long short-term memory (LSTM) network combined with attention mechanism to detect anomalies. Zhang [36] proposed a two-layer parallel learning cross-fusion deep learning model (PCCN), which uses feature fusion technology to improve the extraction of features from small sample data, and experiments on ablation experiments showed good performance. Zhong [37] proposed HELAD, a network anomaly traffic detection algorithm integrating multiple deep learning techniques. Although HELAD has better adaptability and detection accuracy, its bit error rate is slightly higher.

2.2 Transformer architecture

In 2018, transformer architecture [38] first appeared in the field of natural language processing (NLP), and it has occupied an important position in the field of NLP. Transformer architecture has been continuously improved by subsequent scholars [39]. Vaswani [40] first constructed transformer architecture based on attention mechanism. Devlin et al. [41] proposed BERT, a new language representation model, which pretrains a transformer from unmarked text through joint adjustments of left and right contexts. BERT got the latest results from 11 natural language processing tasks at the time.

Influenced by the excellent performance of transformer architecture in NLP task, scholars began to extend transformer architecture to the field of computer vision and achieved good results. Chen et al. [42] constructed a sequence transformer to perform regression prediction of pixels and obtained competitive results in the image classification task. In 2020, Dosovitskiy et al. [43] proposed a vision transformer architecture,

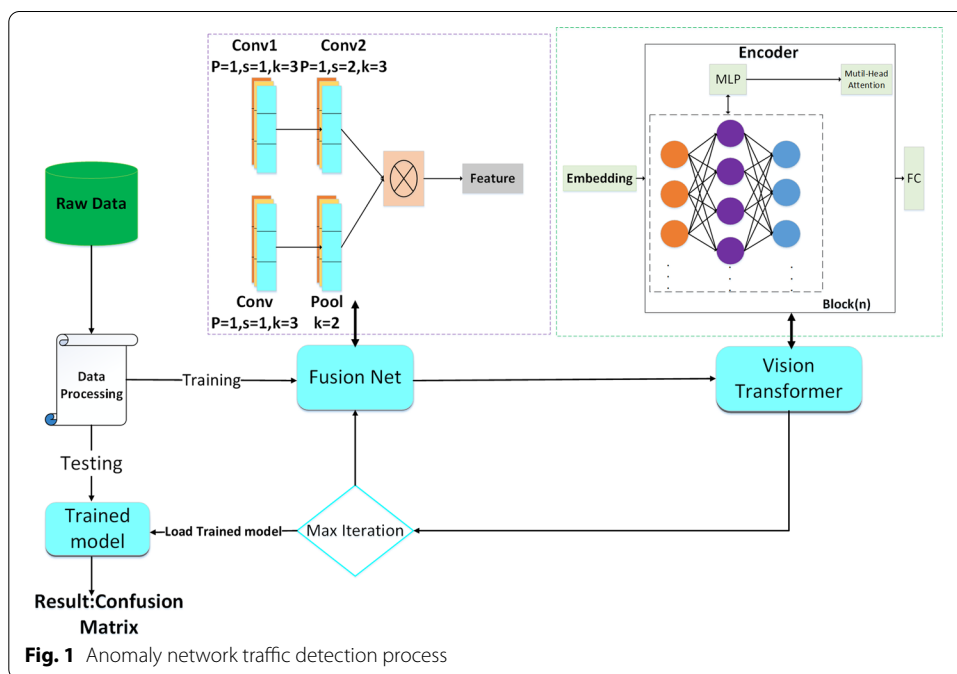
which uses a pure transformer to directly extract the features of image block sequences and obtain the most advanced performance on multiple image recognition reference datasets. Besides the most basic image classification tasks [44], transformer models are gradually applied to various computer vision tasks, and the number of vision models based on transformer architecture has gradually become more and more.

In this paper, the latest intrusion detection model based on feature fusion is improved and integrated into vision transformer architecture, and then a deep learning model (MFVT) that combines feature fusion network with vision transformer architecture is proposed for network anomaly traffic detection. The MFVT takes full advantage of the respective strengths of feature fusion and vision transformer architecture and further improves the detection accuracy of abnormal network traffic by combining with the CPR algorithm proposed by us.

3 Model and methods

This section mainly introduces the CPR algorithm and MFVT model.

In order to improve the processing capacity of existing deep learning models for imbalanced datasets and reduce the required training set resources, in this paper, a new model MFVT and a new raw data processing algorithm CPR were designed. This section mainly introduces the MFVT model and the CPR algorithm. The MFVT model can improve the detection ability of small sample datasets and reduce the training set resources, and the CPR can effectively remove the interference features in the raw data. Figure 1 shows the entire detection process. The MFVT model mainly composed of a feature fusion network and the vision transformer architecture. MFVT can use the raw features of network traffics to automatically learn the differences between different categories of network traffic features to classify anomaly network traffics, but the network model requires that the dimensionality of all input data must be consistent, so an



algorithm named CPR was proposed to extract the raw features of network traffics and intercept the same dimensional data.

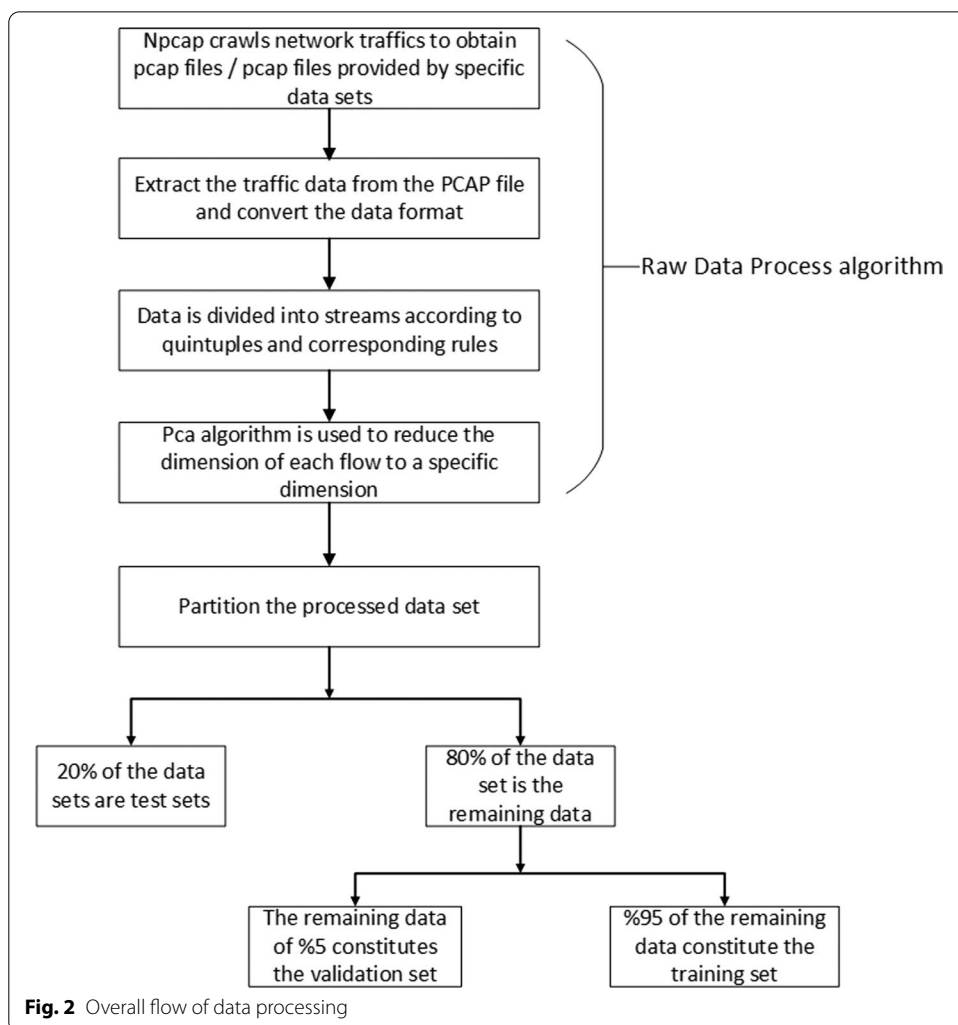
3.1 Data processing

The raw data processing algorithm (CPR) proposed in this paper mainly accomplishes the task of extracting raw traffic data from pcap files and processing them into the two-dimensional matrix that required by the network model [45, 46]. Figure 2 shows the entire data processing process.

Three steps are required to process the raw flow data into a two-dimensional matrix. The specific steps are as follows.

The first step is to extract the raw data of network traffic from the pcap file and then convert the extracted byte type data into binary type data.

In the second step, the converted packets are divided into flows according to the five-tuple, and the number of packets and bytes contained in each packet are limited when dividing the flow. If the number of data packets is insufficient, fill in the preceding item, and if the number of bytes contained in the data packet is insufficient, fill in 0. For the



completion of this step, refer to the paper [47]. Through the above operations, a dataset with fixed dimensions can be obtained. The pseudocode is shown in algorithm 1.

Algorithm 1 Raw data processing

Input: Raw data (pcap);
Output: all_data[];

- 1: **for** each data in pcap **do**
- 2: **if** the same five-tuple could be found in the attack Labels **then**;# Extracting and tagging malicious traffic from pcap files
- 3: Save data and labels into a pcap file
- 4: **end if**
- 5: **end for**
- 6: **for** each data in pcap **do**
- 7: set file name,count=0
- 8: **if** data's five-tuple equal **then**: Save the data into the same flow
- 9: **end if**
- 10: **for** each data in flow **do**: Use wireshark to get data's original hexadecimal data Change the original hexadecimal data into 10hex data Save the 10hex data into mid_data #mid_data is to store each packet data in each flow all_data.append(mid_data) mid_data=[]
- 11: **end for**
- 12: **end for**

In the third step, the network traffic data obtained after the first two steps contain high data dimensions and may have redundant features that are useless for network training, which need to be further extracted. In this paper, the data obtained from the first two steps are directly fed into the PCA algorithm to obtain the data of the required dimensions, and then the data are processed into a two-dimensional matrix. The pseudocode is shown in algorithm 2.

Algorithm 2 Crop and reduce data dimensionality

Input: all_data[];

Output: Data required by neural network;

- 1: data=[] # Store the final processed data
- 2: size # The set number of packets to be intercepted per flow
- 3: length # The maximum number of bytes to be intercepted per packet set
- 4: **for** i=0 to all_data.length **do**
- 5: **if** len(all_data[i]) more than the size **then**;
- 6: Save data and tags into a pcap file;
- 7: **for** j to length **do**
- 8:
- 9: **if** all_data[i][j]==" " **then**; # The number of bytes contained in the packet is less than the number of bytes intercepted fill in 0
- 10: mid_data.append(0)
- 11: **else**
- 12: mid_data.append(all_data[i][j])
- 13: data.append(mid_data)
- 14: **end if**
- 15: Same as above len(all_data[i]) \geq size
- 16: **for** i to (size-len(all_data[i])) **do**
- 17: data.append(The data extracted from the previous)
- 18: **end for**
- 19: **end for**
- 20: **end if**
- 21: **end for**
- 22: data=pca(data, dimension) # Reducing data to a specified dimension using the pca algorithm
- 23: Maxmin_ Normalized(data)
- 24: Save the descended data to the specified csv file

The main idea of PCA is to map the N-dimensional features to the K-dimension, which is a new orthogonal feature, also known as the principal component, and is a

reconstructed K -dimensional feature based on the original N -dimensional features, as shown in Formula 1,2,3,4,5.

$$x_{ij} = x_{ij} - \frac{\sum_{i=1}^n x_{i,j}}{n} \quad 0 < i < n, 0 < j < d \quad (1)$$

$$C = \frac{1}{m} XX^T \quad (2)$$

$$w, b = \text{eig}(c) \quad (3)$$

$$p = \text{select}(\text{sort}(w, b), k) \quad (4)$$

$$Y = PX_{n,d} \quad (5)$$

Formula 1 indicates that the original data X is arranged into a matrix with n rows and D columns, and then the matrix is zero-averaged. x_{ij} represents the data in row i and column j of matrix X . In Formula 2, c represents the covariance matrix of matrix X . Formula 3 expresses getting the eigenvalue and eigenvector of the covariance matrix c , $\text{eig}()$ is the function of getting the eigenvalue and eigenvector, w indicates the obtained eigenvector, and b indicates the corresponding eigenvalue. In Formula 4, the eigenvectors are arranged into a matrix in rows from top to bottom according to the corresponding eigenvalues. The first k rows are taken to form the matrix p , where $\text{sort}()$ is the sorting function and $\text{select}()$ is the selection function. Formula 5 represents the dataset Y obtained after dimension reduction.

3.2 The structure of MFVT

Figure 3 is the overall structure of the MFVT model, which composed of two parts.

First part is the feature fusion network, which is composed of two layers of parallel convolution networks. The first layer is stacked with two convolution layers, the first convolution has a step of 1, the second convolution has a step of 2, and the size of the kernel is 3. The second layer consists of a convolutional layer and a pooling layer, where the convolutional layer has a kernel size of 3 and a step size of 1, and the pooling layer has a step size of 2. The padding size used in the two-layer convolution process is all 1. To make full use of the features extracted by convolution layer and pooling layer, the extracted features are fused to improve the extraction effect of features for small sample data. The whole calculation process of the feature fusion network is shown in Formula 6–16. Formula 6 represents the padding operation, and Formula 7 represents the size change of the output matrix of convolution processing after the padding operation. Under the premise that padding_n is equal to 1, the stride=1 keeps the output size unchanged, and the stride=2 halves the output size.

$$X = \text{Padding}(X_0, 1) = \begin{bmatrix} x_{11} & \cdots & x_{1W} \\ \vdots & \ddots & \vdots \\ x_{H1} & \cdots & x_{HW} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0 & \cdots & 0 \\ x_{11} & \cdots & x_{1W} \\ \vdots & \vdots & \vdots \\ x_{H1} & \cdots & x_{HW} \\ 0 & \cdots & 0 \end{bmatrix} \tag{6}$$

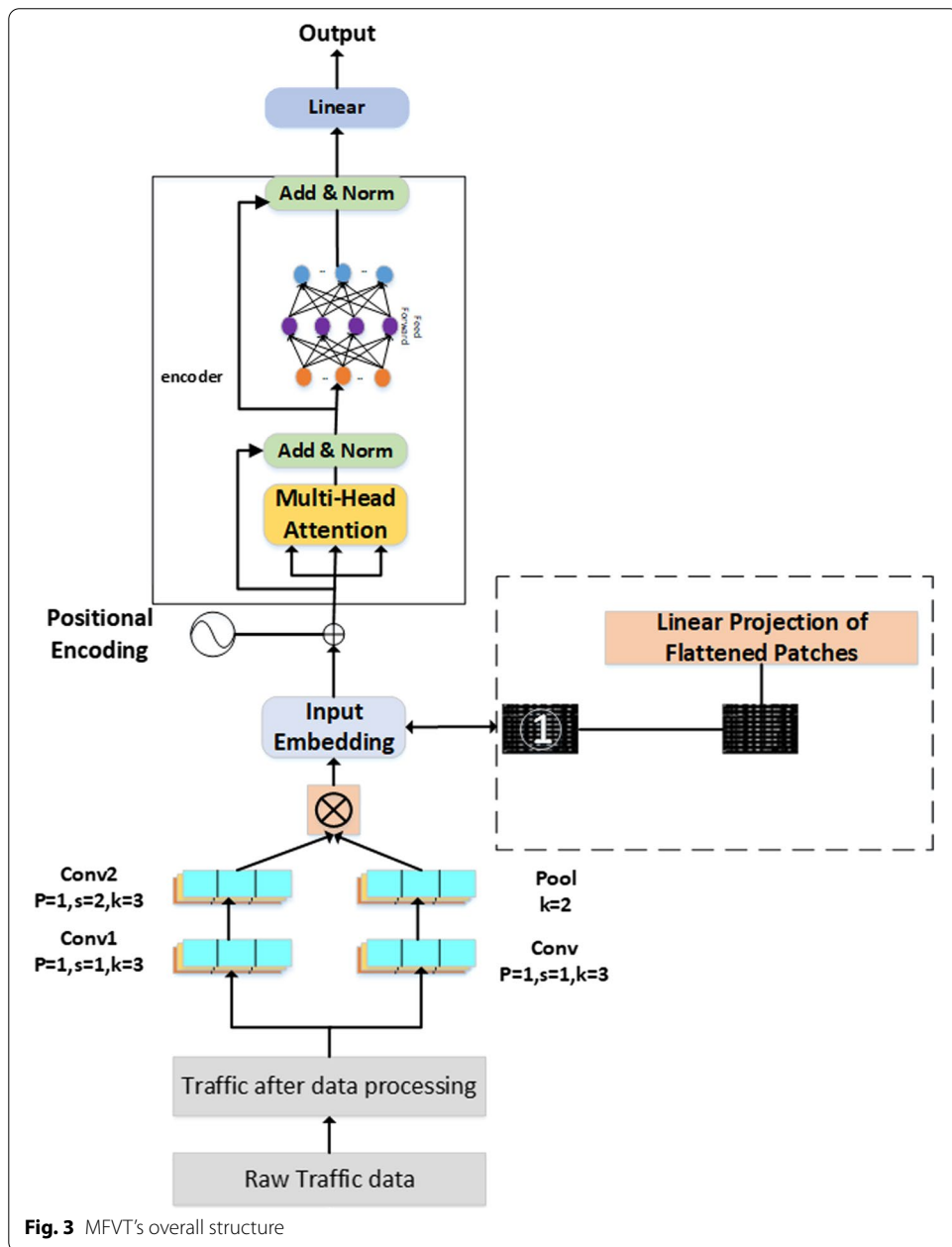


Fig. 3 MFVT's overall structure

$$\begin{aligned}
 H &= \left\lceil \frac{H + \text{padding}_n - w + 1}{\text{stride}} \right\rceil \\
 W &= \left\lceil \frac{W + \text{padding}_n - w + 1}{\text{stride}} \right\rceil
 \end{aligned}
 \tag{7}$$

Where X_O represents the matrix data obtained after the original flow data is processed by the CPR algorithm, because the convolution manipulation will change the size of the input matrix, in order to keep the matrix size unchanged it is necessary to perform the padding operation by Formula 7. X represents the matrix after the padding operation, X_{ij} represents the specific data value in the matrix. W is the width of the matrix, and H is the height.

Formulas 6, 8, 9, 10 represent the entire calculation process of the first layer in the feature fusion network. Formulas 8 and 10 represent the convolution operation, V represents the convolution kernel matrix, v_{ij} represents the specific value in the convolution kernel matrix, and k represents the kernel sizes. X_1^1 represents the eigenmatrix obtained after the first convolution operation. Since the stride in Formula 7 is 1, the output size remains unchanged. X_1^2 represents the matrix obtained after the padding operation of X_1^1 , and X_1^3 represents the eigenmatrix obtained after the second convolution, and the output size is halved because the stride in Formula 7 is 2.

$$X_1^1 = X \odot V = \begin{bmatrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{k1} & \cdots & x_{kk} \end{bmatrix} \odot \begin{bmatrix} v_{11} & \cdots & v_{1k} \\ \vdots & \ddots & \vdots \\ v_{k1} & \cdots & v_{kk} \end{bmatrix}
 \tag{8}$$

$$X_1^2 = \text{padding} \left(X_1^1, 1 \right)
 \tag{9}$$

$$X_1^3 = X_1^2 \odot V = \begin{bmatrix} x_{11}^2 & \cdots & x_{1k}^2 \\ \vdots & \ddots & \vdots \\ x_{k1}^2 & \cdots & x_{kk}^2 \end{bmatrix} \odot \begin{bmatrix} v_{11} & \cdots & v_{1k} \\ \vdots & \ddots & \vdots \\ v_{k1} & \cdots & v_{kk} \end{bmatrix}
 \tag{10}$$

Formulas 6, 11, 12 represent the entire computational process of the second layer in the feature fusion network, where X_2^1 denotes the feature matrix extracted after the convolution operation, the stride=1 does not change the output size, and X_2^2 denotes the feature matrix obtained after the maximum pooling operation, which halves the size of the output feature matrix.

$$X_2^1 = X \odot V = \begin{bmatrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{k1} & \cdots & x_{kk} \end{bmatrix} \odot \begin{bmatrix} v_{11} & \cdots & v_{1k} \\ \vdots & \ddots & \vdots \\ v_{k1} & \cdots & v_{kk} \end{bmatrix}
 \tag{11}$$

$$X_2^2 = \text{Maxpooling} \left(X_2^1 \right) = \frac{\max \left\{ x_{ij}^1 \right\}}{i, j \in [1, k]}
 \tag{12}$$

Formula 13 shows the scale changes of the features extracted from the first and second layers of the feature fusion network. Formula 14 represents the specific process of fusing

the first layer with the second layer features. The fusion refers to the summation of the number of channels, but the data must be kept consistent except for the number of channels. C represents the number of channels, $C(1)$ represents the number of channels is 1, $C(32)$ represents the number of channels is 32 and so on, X_f represents the features extracted by the feature fusion network.

$$(C(1), H, W) \Rightarrow \left(C(32), \frac{H}{2}, \frac{W}{2} \right) \tag{13}$$

$$\begin{aligned} X_f &= \left(C(32), \frac{H}{2}, \frac{W}{2} \right) \oplus \left(C(32), \frac{H}{2}, \frac{W}{2} \right) \\ &= \left(C(32 + 32), \frac{H}{2}, \frac{W}{2} \right) \end{aligned} \tag{14}$$

The second part is composed of the vision transformer architecture. To combine vision transformer architecture with feature fusion network, the structure of vision transformer is modified in this paper. The main methods used include feature embedding, learnable embedding, and transformer encoder.

For feature embedding, standard transformer accepts sequence of token embeddings as input. To process the feature X_f learned by the feature fusion network, we reconstructed X_f into a flattened 2D block sequence X_p . Formula 15 is a specific variation of the formula, the same as NLP, will be added to the sequence of images in the token classification, the sequence of images is cut into multiple patches by a picture to get the number of patches where p indicates.

$$\begin{aligned} X_f &\in R^{C(64) \times \frac{H}{2} \times \frac{W}{2}} \\ X_p &\in R^{N \times (P^2 C(64))} \\ N &= \frac{\frac{H}{2} \times \frac{W}{2}}{p^2} \\ X_f &\rightarrow X_p \end{aligned} \tag{15}$$

Learnable embedding, a learnable embedding $z_0^0 = x_{class}$ is preset for the feature block embedding sequence, x_{class} denotes the category vector whose state/feature Z_L^0 at the transformer encoder output is used as the feature representation y , as shown in Formula 21. Learnable embedding is randomly initialized at the beginning of training and obtained by training.

Transformer encoder consists of several blocks, each containing a multi-head attention block and a multi-layer perceptron (MLP) block, with normalization applied before each block and residual concatenation applied after each block. Figure 4 shows the structure of head attention. Formula 16, 17 show how to get the multi-head attention values by head attention. Where W_i^Q, W_i^K, W_i^V and W^O are all weight matrices.

$$\text{Head}_i = \text{Attention} \left(QW_i^Q, KW_i^K, VW_i^V \right) \tag{16}$$

$$\text{Multihead} (Q, K, V) = \text{Concat} (\text{head}_1, \dots, \text{head}_b) W^O \tag{17}$$

Finally, the embedding vectors that combine category vectors and feature block embedding can be input into transformer encoder. The encoder built up by blocks can extract data features for classification just like CNN. The whole calculation process is shown in Formula 18, 19, 20, 21.

$$Z_0 = [X_{class}; X_p^1 E], \quad E \in R^{(P^2 \cdot C) \times D}, E_{pos} \in R^{(N+1) \times D} \tag{18}$$

$$Z'_l = MSA(LN(Z_{l-})) + Z_{l-1}, \quad l = 1 \dots L \tag{19}$$

$$Z_l = MLP(LN(Z'_l)) + Z'_l, \quad l = 1 \dots L \tag{20}$$

$$y = LN(Z_L^0) \tag{21}$$

The feature embedding block $X_p^1 E$ and the category vector X_{class} form the embedding input vector Z_0 . Formula 19 adopts skip connection, where MAS represents multi-head attention operation, LN represents normalization operation, L represents repeatable times, and Z'_l represents the l th output. Formula 20 adopts skip connection, MLP represents the multi-layer perceptron block, L represents repeatable times, and Z_l represents the l th output. y represents the feature representation.

4 Experiments and results analysis

This section first introduces the experimental environment, the datasets IDS 2017 and IDS 2012 used in the experiments, the evaluation criteria used in the experiments and finally specifies the ablation experiments and some details of the experiments. In the ablation experiments, a series of advanced models were compared with the MFVT model.

4.1 The experimental environment of this paper

In this paper, ablation experiments were conducted on the MFVT model and CPR data processing algorithm under the environment shown in Table 2.

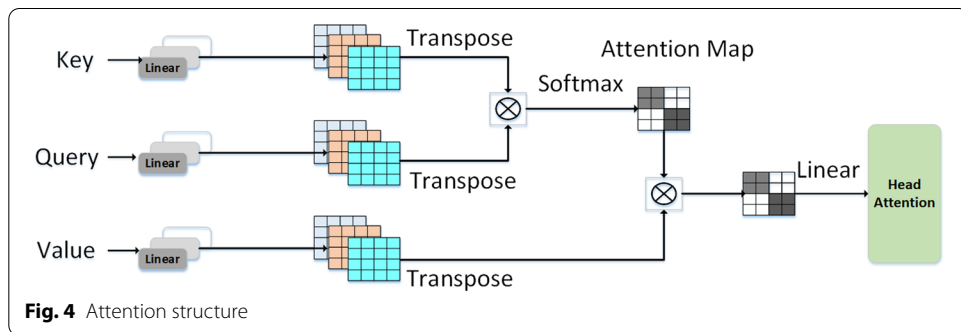


Fig. 4 Attention structure

4.2 Datasets

In this paper, A series of ablation experiments were designed using both IDS 2012 and IDS 2017 datasets.

The IDS 2012 dataset contains a week of network activity including both normal and malicious activity, with three days consisting of all normal traffics and the remaining four days consisting of a large amount of normal traffics with a specific type of attack traffics. IDS 2012 dataset contains attack traffic including internal penetration, HTTP denial of service, distributed denial of service using IRC botnet, and brute force cracking of SSH [18].

The IDS 2017 data collection period lasts for five days from 9am on Monday, July 3, 2017, to 5pm on Friday, July 7, 2017, of which Mondays only include normal traffic. The attacks implemented included brute force FTP, brute force SSH, DoS, Heartbleed, web attack, infiltration, Botnet, and DDoS [17].

Figure 5a is a bar chart of the amount of various attack traffics contained in the IDS 2017 dataset, and Fig. 5b is a pie chart of the amount of various attack traffics contained in the IDS 2012 dataset. It is observed from the figures that both IDS 2017 and IDS 2012 datasets have serious data imbalance problems. The data volume of DDOS, Hulk, and PortScan attacks in the IDS 2017 dataset is significantly larger than that of other types of attacks. The data volume of Infiltrator attacks in the IDS 2012 dataset directly accounts for 55% of the dataset

4.3 Evaluation metrics

Authoritative evaluation metrics must be used to judge the merits of a network anomaly traffic detection method. The effectiveness of the machine learning-based network anomaly traffic detection algorithm can be evaluated by the metrics shown in Formula 26, 25, 23, 22, 24 [48]. TP represents the positive sample predicted to be positive by the model, which can be called the accuracy rate judged to be true. TN represents the negative sample predicted to be negative by the model, which can be referred to as the percentage of correct judgments that are false. FP represents the negative sample predicted by the model to be positive, which can be referred to as the false alarm rate. FN represents the positive sample predicted to be negative by the model, which can be referred to as the underreporting rate [37].

$$\text{Recall} = \frac{TP}{TP + FN} \quad (22)$$

Table 2 Experimental environment of this paper

CPU:	i7-10875H CPU@2.30GHz 2.30GHz
RAM:	16G
GPU:	RTX 2060 6G
Compiler environment:	Python 3.8.2
OS:	Windows 10

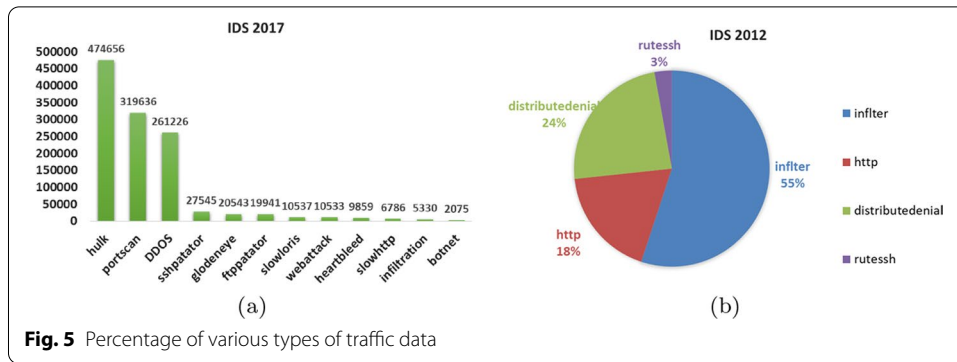


Fig. 5 Percentage of various types of traffic data

$$\text{Precision} = \frac{TP}{TP + FP} \tag{23}$$

$$F1 - \text{Measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \tag{24}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{25}$$

$$FPR = \frac{FP}{FP + TN} \tag{26}$$

4.4 Ablation experiment and results analysis

In this paper, two datasets of IDS 2012 and IDS 2017 were used for ablation experiments. In addition, this paper also carried out an exploratory study on the impact of model optimization methods on MFVT model detection performance on IDS 2012 dataset.

In the MFVT model, the size of the kernels used in the convolutional neural network is 3 * 3, the segmentation size set in the vision transformer architecture is 11 * 11, the number of the head in the multi-head attention is 12, and the number of blocks in the encoder is 12. In the process of model training, the data input batch used in this paper is 256, the epoch of the training iteration is set to 100, and the stochastic gradient descent (SGD) optimizer is used to accelerate the network convergence. The momentum is fixed at 0.9, the learning rate is fixed at 3e-2, weight_decay Set to 0, the loss function uses CrossEntropyLoss. All ablation experiments and results will be described in detail below.

4.4.1 Ablation experiment based on IDS 2012

Figure 6 shows the parameter changes of MFVT model when using IDS 2012 dataset for training, including training loss, verification loss, and verification accuracy. As shown in the picture, the convergence speed of MFVT model is fast, but there are large fluctuations in the later stages of training.

Table 3 shows the experimental results based on IDS 2012 dataset. It is obvious from the table that the MFVT model combined with CPR algorithm proposed in this paper is superior to other methods on all evaluation metrics, reaching the state-of-the-art level.

It can also be concluded from the table that MFVT model has superior performance, and its detection accuracy is only slightly worse than that of DT (Decision Tree), but it has higher precision. To better demonstrate the ability of the MFVT model to deal with imbalanced data, the experimental results of all evaluation metrics of the MFVT model in various types of attack traffic are shown in Table 4.

Combining the (B) in Fig. 5 and Table 4 (the experimental results of Infiltrating and Distributed denial, which account for a relatively large proportion, have been marked in bold), it can be concluded that the traffic of HTTP and rutesh, which account for a relatively small proportion, still obtains good experimental results. It shows that MFVT model has strong ability to recognize small sample data. The detection performance is further improved by combining the CPR algorithm with the MFVT model.

4.4.2 Ablation experiment based on IDS 2017

The IDS 2012 dataset contains fewer types of attack traffic, and the effectiveness of the MFVT model and the data processing algorithm CPR are demonstrated to be not generalizable on this dataset only. So, ablation experiments also were performed on the more complex IDS 2017 dataset. Figures 7 and 8 are the results of the ablation experiment, from which it can be seen that the accuracy, recall, F1-score and accuracy of the MFVT model and the combination of MFVT model and CPR algorithm all reached nearly 100%, which was significantly better than other comparison models. Figure 7 shows that the detection results obtained by MFVT model and the combination of MFVT model and CPR algorithm are close to 100% in the evaluation criteria, which is significantly better than other comparison models. The comparison of the FPR between the MFVT model and other comparison models is shown in Fig. 8, from which it can be seen that the MFVT model is still the best. Combined with Fig. 5a and Table 5 (experimental results of DDos, Hulk and Portscan, which account for a large proportion of attacks, have been marked in bold), it can be concluded that the MFVT model combined with the CPR algorithm has a better ability to recognize small samples.

To further demonstrate the error of the prediction results of the MFVT model proposed in this paper combined with CPR, the experimental results were made into the heat map shown in Fig. 9. From the heat map, the performance of the MFVT model combined with CPR is very high, and the prediction error rate is extremely low.

To verify that the MFVT model can reduce the sample resources required for training, we tested it on the IDS 2017 dataset by reducing the training set data volume according to Formula 27 with all other conditions held constant, $data_0$ is the initial assigned training set data volume, $data_n$ is the updated data volume, and n is taken according to Formula 28, where n_0 is the initial value of n equal to 0.9, and N takes values in the range of 1–7. Table 6 shows the test results.

$$Data_n = (1 - 0.1 * n) * data_0 \quad (27)$$

$$n = n_0 - 0.1N \quad (28)$$

As can be seen from Fig. 10, when the training set data amount is reduced to 80% of the original training set data amount, the impact on the overall accuracy of the test set is very small. Through this experiment, it is proved that the MFVT model combined with

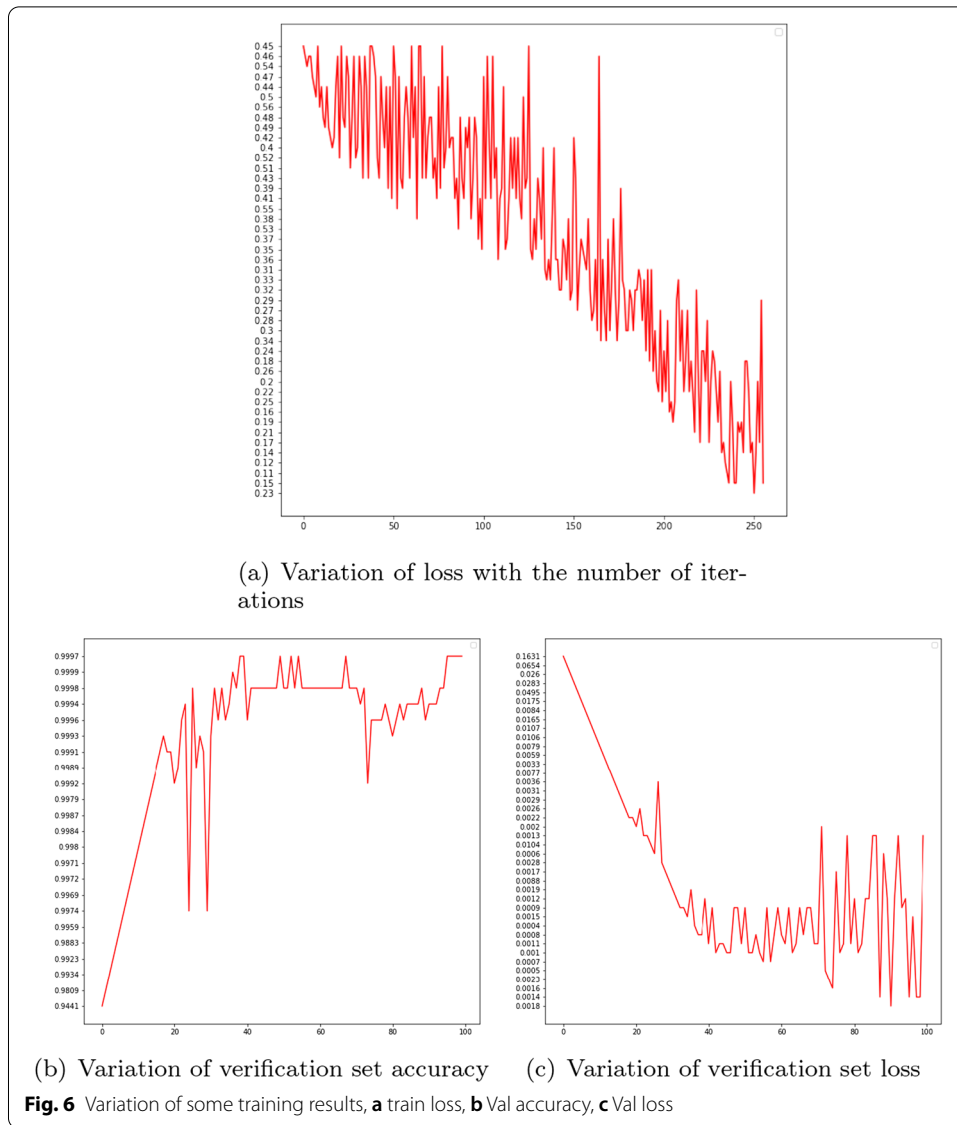
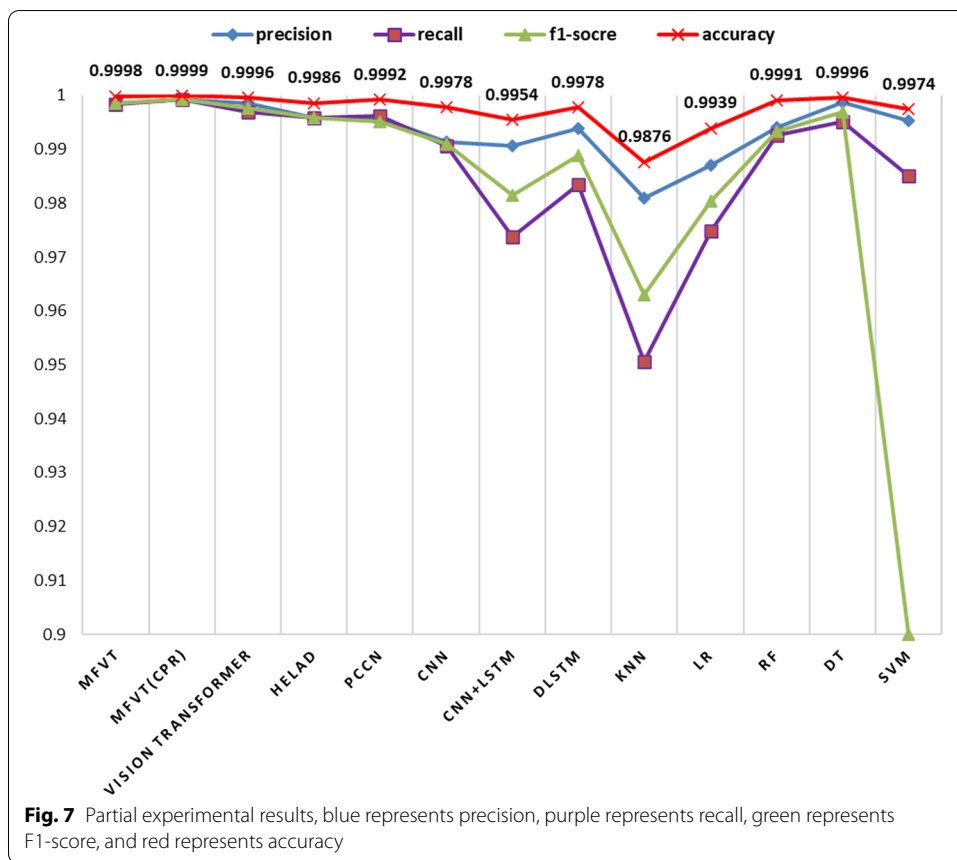


Table 3 Experimental results on the 2012 dataset

Methods	Precision	Recall	F1-score	FPR	Accuracy
MFVT	0.9986	0.9975	0.998	0.000525	0.9988
MFVT (CPR)	0.9995	0.9994	0.9995	0.000175	0.9996
vision transformer	0.9984	0.9977	0.998	0.000625	0.9985
PCCN	0.9987	0.9979	0.9983	0.000575	0.9986
CNN	0.9958	0.9942	0.9949	0.00145	0.9962
CNN_LSTM	0.9949	0.9936	0.9942	0.001775	0.9951
DLSTM	0.9939	0.9928	0.9933	0.00195	0.9944
KNN	0.993	0.9903	0.9917	0.002125	0.9939
LR	0.9891	0.9902	0.9897	0.00315	0.9909
RF	0.9973	0.9966	0.9969	0.00085	0.9979
DT	0.9984	0.9984	0.9984	0.000375	0.999
SVM	0.9943	0.9937	0.994	0.0018	0.9949

Table 4 Performance of MFVT model and CPR algorithm in each category in IDS 2012 data

Methods	Class	Precision	Recall	F1-socre	False alarm rate
MFVT	Infiltrating	0.999	0.9994	0.9992	0.0012
	http	0.9984	0.9968	0.9976	0.0004
	Distributeddenial	0.9985	0.9992	0.9988	0.0005
	rutessh	0.9984	0.9945	0.9965	0
MFVT(CPR)	Infiltrating	0.9995	0.9998	0.9997	0.0006
	http	0.9995	0.9985	0.999	0.0001
	Distributeddenial	0.9999	0.9999	0.9999	0
	rutessh	0.9992	0.9992	0.9992	0



CPR algorithm can effectively reduce the training resources and maintain the accuracy of the test set as much as possible.

4.4.3 Optimization of MFVT model

In the conclusion of this section, it is hoped to further improve the detection accuracy and the stability of the model by increasing the training epochs and continuously adjusting the learning rate (*lr*) during the training process. Thus, IDS 2012 is used as the ablation experiment dataset, which takes less time to train than IDS 2017. Two sets of

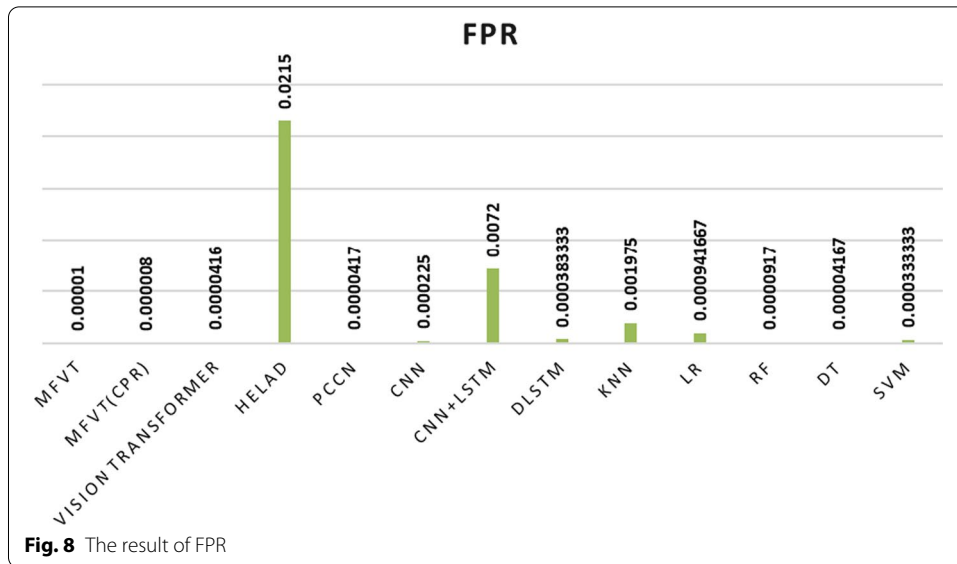


Table 5 Performance of MFVT(CPR) in the IDS 2017 dataset

Class	Precision	Recall	F1-score
botnet	0.9928	1	0.9964
DDoS	0.9999	0.9999	0.9999
goldeneye	0.9995	0.9998	0.9996
hulk	0.9999	1	0.9999
slowhttp	0.9993	0.9971	0.9982
slowloris	1	0.999	0.9995
ftppatator	1	0.9992	0.9996
heartbleed	1	0.9995	0.9997
infiltration	1	0.9981	0.9991
portscan	1	1	1
sshpatorator	0.9996	1	0.9998
webattack	0.9991	0.9991	0.9991

experiments were conducted. In the first group, our model was trained 1,000 times and the results were recorded every 100 times.

In the second group, based on the first group, *lr* is changed 100 times per iteration according to Formula 29, where lr_i is the learning rate changed every time according to the formula, lr_0 is the initial learning rate, and the epoch is every hundred iterations. To ensure the rigor of the experiment, the values were obtained after conducting the two sets of experiments several times. It can be seen from Fig. 11 that both the increase of training epochs and the change of *lr* can get better prediction accuracy in some intermediate results, but the experimental results tend to be stable in the end. In comparison, the variation of *lr* will make the variation of experimental results more stable.

$$lr_i = 0.95^{epoch/10} \cdot lr_0 \tag{29}$$

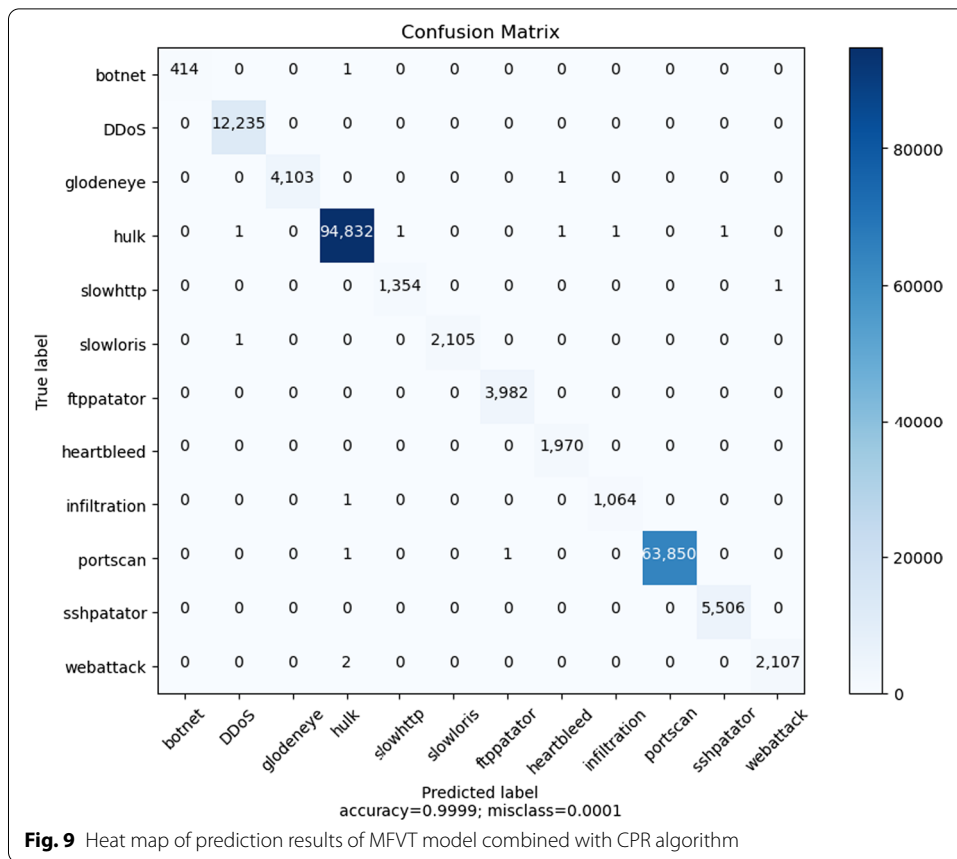


Fig. 9 Heat map of prediction results of MFVT model combined with CPR algorithm

Table 6 Test results-% of original training data

Methods	100%	90%	80%	70%	60%	50%	40%	30%
MFVT(CPR)	0.999932	0.999931	0.999929	0.999845	0.999839	0.999783	0.999532	0.999321

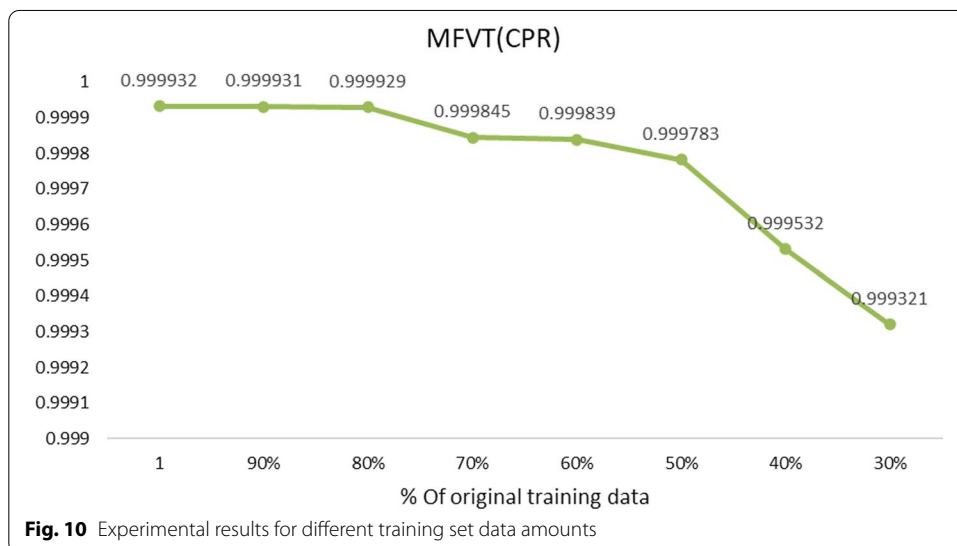
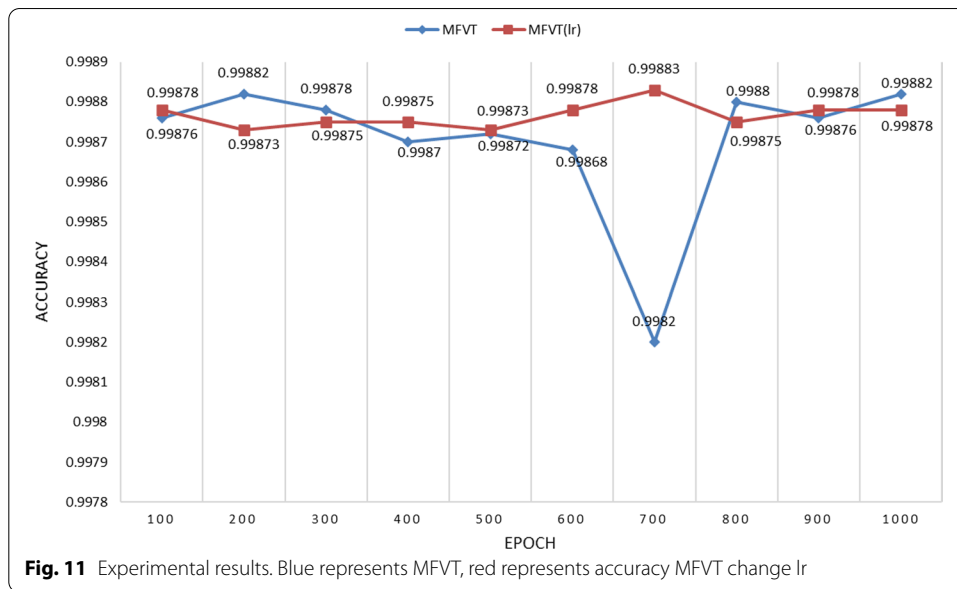


Fig. 10 Experimental results for different training set data amounts



5 Results and discussion

Since most of the deep learning models need a lot of training resources, a network anomaly traffic detection model (MFVT) which combining a feature fusion network with the vision transformer architecture was proposed. MFVT can reduce training resources while maintaining high detection accuracy. In this paper, a new raw traffic data extraction algorithm (CRP) was proposed. The MFVT model combined with the CRP algorithm achieved nearly 100% detection accuracy on both datasets IDS 2012 and IDS 2017, and with much better performance than the other methods in the comparison experiments. The MFVT model combined with the CRP algorithm is more capable of handling imbalanced datasets and can further improves the detection accuracy of the experiment.

Although the MFVT model combined with the CRP algorithm has an excellent performance in the field of anomaly traffic detection, the scalability of the model is weak and the detection accuracy of new types of attack traffic that do not appear in the training set needs to be improved in the face of the increasingly complex network environment and the emergence of new attack types.

Considering the importance and practical significance of scalability, the scalability of the MFVT model will be further improved in the future to enhance the practical value and practical significance of the model.

Abbreviations

MFVT: An anomaly traffic detection method merging feature fusion network and vision transformer architecture; CPR: A new raw traffic features extraction method; PCA: Principal component analysis.

Authors' contributions

All authors read and participated in the manuscript's completion. All authors read and approved the final manuscript.

Funding

This research is supported by the National Natural Science Foundation of China under Grants 61873160, 61672338 and Natural Science Foundation of Shanghai under Grant 21ZR1426500.

Declarations

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Engineering, Shanghai Maritime University, Shanghai, China. ²Department of Information Engineering and Computer Science, Shanghai Maritime University, Taichung, Taiwan.

Received: 5 September 2021 Accepted: 10 March 2022

Published online: 25 April 2022

References

1. D. Han, N. Pan, K.-C. Li, A traceable and revocable ciphertext-policy attribute-based encryption scheme based on privacy protection. *IEEE Trans. Depend. Secure Comput.* (2020)
2. M. Cui, D. Han, J. Wang, An efficient and safe road condition monitoring authentication scheme based on fog computing. *IEEE Internet Things J.* **6**(5), 9076–9084 (2019)
3. Q. Tian, D. Han, K.-C. Li, X. Liu, L. Duan, A. Castiglione, An intrusion detection approach based on improved deep belief network. *Appl. Intell.* **50**(10), 3162–3178 (2020)
4. L. Hung-Jen, R.L. Chun-Hung, L. Ying-Chih, T. Kuang-Yuan, Intrusion detection system: a comprehensive review. *J. Netw. Comput. Appl.* **36**, 16–24 (2013)
5. D. Li, D. Han, Z. Zheng, T.-H. Weng, H. Li, H. Liu, A. Castiglione, K.-C. Li, Moocschain: A blockchain-based secure storage and sharing scheme for moocs learning. *Comput. Stand. Interfaces*, 103597 (2021)
6. D.J. Weller-Fahy, B.J. Borghetti, A.A. Sodemann, A survey of distance and similarity measures used within network intrusion anomaly detection. *IEEE Commun. Surv. Tutor.* **17**, 70–91 (2015)
7. A. Abraham, C. Grosan, C. Martin-Vide, Evolutionary design of intrusion detection programs. *Int. J. Netw. Secur.* **4**, 328–339 (2007)
8. S. Anwar, J. Mohamad Zain, M. Zolkipli, Z. Inayat, S. Khan, B. Anthony Jnr, V. Chang, From intrusion detection to an intrusion response system: Fundamentals, requirements, and future directions. *Algorithms* **10**, 39 (2017)
9. W. Zhang, D. Han, K.-C. Li, F.I. Massetto, Wireless sensor network intrusion detection system based on mk-elm. *Soft Computing*, 1–14 (2020)
10. W. Liang, L. Xiao, K. Zhang, M. Tang, D. He, K.-C. Li, Data fusion approach for collaborative anomaly intrusion detection in blockchain-based systems. *IEEE Internet of Things J.* (2021)
11. A. Ajith, G. Crina, M.V. Carlos, A survey of network anomaly detection techniques. *J. Netw. Comput. Appl.* **60**, 19–31 (2016)
12. J. Zhang, C. Chao, X. Yang, W. Zhou, X. Yong, Internet traffic classification by aggregating correlated naive bayes predictions. *IEEE Trans. Inf. Forens. Secur.* **8**, 5–15 (2013)
13. Y. Zhang, X. Chen, L. Jin, X. Wang, D. Guo, Network intrusion detection: based on deep hierarchical network and original flow data. *IEEE Access* **7**, 37004–37016 (2019)
14. H. Liu, D. Han, D. Li, Behavior analysis and blockchain based trust management in vanets. *J. Parallel Distrib. Comput.* **151**, 61–69 (2021)
15. K. Oksuz, B.C. Cam, S. Kalkan, E. Akbas, Imbalance problems in object detection: a review. *IEEE Trans. Pattern Anal. Mach. Intell.* pp. 1–1 (2020)
16. A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, N. Houlsby, An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
17. I. Sharafaldin, A.H. Lashkari, A.A. Ghorbani, Toward generating a new intrusion detection dataset and intrusion traffic characterization. *ICISp* **1**, 108–116 (2018)
18. A. Shiravi, H. Shiravi, M. Tavallaee, A.A. Ghorbani, Toward developing a systematic approach to generate benchmark datasets for intrusion detection. *Comput. Secur.* **31**(3), 357–374 (2012)
19. L.I. Smith, A tutorial on principal components analysis. *Inf. Fusion* **51**, 52 (2002)
20. D. Han, Y. Zhu, D. Li, W. Liang, A. Soury, K.-C. Li, A blockchain-based auditable access control system for private data in service-centric iot environments. *IEEE Trans. Ind. Inform.* (2021)
21. W. Liang, Z. Ning, S. Xie, Y. Hu, S. Lu, D. Zhang, Secure fusion approach for the internet of things in smart autonomous multi-robot systems. *Inf. Sci.* **579**, 468–482 (2021)
22. H. Li, D. Han, M. Tang, A privacy-preserving storage scheme for logistics data with assistance of blockchain. *IEEE Internet of Things J.* (2021)
23. X. Chen, W. Liang, J. Xu, C. Wang, K.-C. Li, M. Qiu, An efficient service recommendation algorithm for cyber-physical-social systems. *IEEE Trans. Netw. Sci. Eng.* (2021)
24. J.P. Anderson, *Computer security threat monitoring and surveillance* (1980)
25. C.L. Yin, Y.F. Zhu, J.L. Fei, X.Z. He, A deep learning approach for intrusion detection using recurrent neural networks. *IEEE Access*, pp. 1–1 (2017)
26. F. Kuang, W. Xu, S. Zhang, A novel hybrid kpca and svm with ga model for intrusion detection. *Appl. Soft Comput.* **18**, 178–184 (2014)
27. R.R. Reddy, Y. Ramadevi, K. Sunitha, Effective discriminant function for intrusion detection using svm. In: 2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI), pp. 1148–1153 (2016)
28. W. Li, P. Yi, Y. Wu, L. Pan, J. Li, A new intrusion detection system based on knn classification algorithm in wireless sensor network. *J. Electr. Comput. Eng.* **2014** (2014)
29. N. Farnaaz, M.A. Jabbar, Random forest modeling for network intrusion detection system. *Procedia Comput. Sci.* **89**, 213–217 (2016)

30. Random-forests-based network intrusion detection systems, *IEEE Trans. Syst. Man Cybernet. Part C* **38**, 649–659 (2008)
31. Y. Dhote, S. Agrawal, A.J. Deen, A survey on feature selection techniques for internet traffic classification. In: *International Conference on Computational Intelligence & Communication Networks*, pp. 1375–1380 (2015). IEEE
32. H. Zhang, G. Lu, M.T. Qassrawi, Y. Zhang, X. Yu, Feature selection for optimizing traffic classification. *Comput. Commun.* **35**, 1457–1471 (2012)
33. A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks. *Adv. Neural Inf. Process. Syst.* **25**, 1097–1105 (2012)
34. Q. Yan, M. Wang, W. Huang, X. Luo, F.R. Yu, Automatically synthesizing dos attack traces using generative adversarial networks. *Int. J. Mach. Learn. Cybern.* **10**, 3387–3396 (2019)
35. P. Lin, K. Ye, C.-Z. Xu, Dynamic network anomaly detection system by using deep learning techniques. In: *International Conference on Cloud Computing*, pp. 161–176 (2019). Springer
36. Y. Zhang, X. Chen, D. Guo, M. Song, X. Wang, Pccn: Parallel cross convolutional neural network for abnormal network traffic flows detection in multi-class imbalanced network traffic flows. *IEEE Access*, pp. 1–1 (2019)
37. Y. Zhong, W. Chen, Z. Wang, Y. Chen, K. Li, Helad: A novel network anomaly detection model based on heterogeneous ensemble learning. *Comput. Netw.* **169**, 107049 (2019)
38. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need. In: *Advances in Neural Information Processing Systems*, pp. 5998–6008 (2017)
39. K. Han, Y. Wang, H. Chen, X. Chen, D. Tao, A survey on visual transformer. *arXiv preprint arXiv:2012.12556* (2020)
40. A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever et al., Language models are unsupervised multitask learners. *OpenAI blog* **1**, 9 (2019)
41. M. Kim, G. Kim, S.-W. Lee, J.-W. Ha, St-bert: Cross-modal language model pre-training for end-to-end spoken language understanding. In: *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7478–7482 (2021). IEEE
42. Y. Chang, Z. Huang, Q. Shen, The same size dilated attention network for keypoint detection. In: *International Conference on Artificial Neural Networks*, pp. 471–483 (2019). Springer
43. J. Chung, C. Gulcehre, K. Cho, Y. Bengio, Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014)
44. W. Liang, J. Long, K.-C. Li, J. Xu, N. Ma, X. Lei, A fast defogging image recognition algorithm based on bilateral hybrid filtering. *ACM transactions on multimedia computing, communications, and applications (TOMM)* **17**, 1–16 (2021)
45. T. Xiao, D. Han, J. He, K.-C. Li, R.F. de Mello, Multi-keyword ranked search based on mapping set matching in cloud ciphertext storage system. *Connect. Sci.* **33**, 95–112 (2021)
46. W. Liang, D. Zhang, X. Lei, M. Tang, K.-C. Li, A. Zomaya, Circuit copyright blockchain: Blockchain-based homomorphic encryption for ip circuit protection. *IEEE Trans. Emerg. Top. Comput.* (2020)
47. M. Li, D. Han, X. Yin, H. Liu, D. Li: Design and implementation of an anomaly network traffic detection model integrating temporal and spatial features. *Secur. Commun. Netw.* **2021** (2021)
48. M. Cui, D. Han, J. Wang, K.-C. Li, C.-C. Chang, Arfv: an efficient shared data auditing scheme supporting revocation for fog-assisted vehicular ad-hoc networks. *IEEE Trans. Veh. Technol.* **69**(12), 15815–15827 (2020)

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)
