

# Detecting Unknown HTTP-based Malicious Communication Behavior via Generated Adversarial Flows and Hierarchical Traffic Features

Xiaochun Yun<sup>a</sup>, Jiang Xie<sup>b,d</sup>, Shuhao Li<sup>b,c,d,\*</sup>, Yongzheng Zhang<sup>b,c,d</sup>, Peishuai Sun<sup>b,d</sup>

<sup>a</sup>National Computer Network Emergency Response Technical Team/Coordination Center of China, Beijing, China

<sup>b</sup>Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

<sup>c</sup>Key Laboratory of Network Assessment Technology, University of Chinese Academy of Sciences, Beijing, China

<sup>d</sup>School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

---

## Abstract

Malicious communication behavior is the network communication behavior generated by malware (botnet, spyware, *etc.*) after victim devices are infected. Experienced adversaries often hide malicious information in HTTP traffic to evade detection. However, related detection methods have inadequate generalization ability because they are usually based on artificial feature engineering and outmoded datasets. In this paper, we propose an HTTP-based Malicious Communication traffic Detection Model (HMCD-Model) based on generated adversarial flows and hierarchical traffic features. HMCD-Model consists of two parts. The first is a generation algorithm based on WGAN-GP to generate HTTP-based malicious communication traffic for data enhancement. The second is a hybrid neural network based on CNN and LSTM to extract hierarchical spatial-temporal features of HTTP-based traffic. In addition, we collect and publish a dataset, HMCT-2020, which consists of large-scale malicious and benign traffic during three years (2018-2020). Taking the data in HMCT-2020(18) as the training set and the data in other datasets as the test set, the experimental results show that the HMCD-Model can effectively detect unknown HTTP-based malicious communication traffic. It can reach  $F1 \approx 98.66\%$  in the dataset HMCT-2020(19-20),  $F1 \approx 90.69\%$  in the public dataset CIC-IDS-2017 and  $F1 \approx 83.66\%$  in the real traffic, which is 20+% higher than other representative methods on average. This validates that HMCD-Model has the ability to discover unknown HTTP-based malicious communication behavior.

*Keywords:* Malicious Behavior Detection, CNN, LSTM, GAN, Hierarchical Features

---

\*The corresponding author of this paper is Shuhao Li.

*Email addresses:* [yunxiaochun@cert.org.cn](mailto:yunxiaochun@cert.org.cn) (Xiaochun Yun), [xiexiang@iie.ac.cn](mailto:xiexiang@iie.ac.cn) (Jiang Xie), [lishuhao@iie.ac.cn](mailto:lishuhao@iie.ac.cn) (Shuhao Li), [zhangyongzheng@iie.ac.cn](mailto:zhangyongzheng@iie.ac.cn) (Yongzheng Zhang), [sunpeishuai@iie.ac.cn](mailto:sunpeishuai@iie.ac.cn) (Peishuai Sun)

## 1. Introduction

Malicious traffic detection generated by malware is one of the hot issues in cyber security [15, 22]. In this paper, we call the network communication behavior generated by malware (botnet, spyware, *etc.*) after malware infects the victim device as *malicious communication behavior*. One of the main carriers of these behavior is HTTP traffic[44]. Experienced adversaries construct HTTP-based malicious communication traffic by imitating the network behavior of benign users and hiding malicious information into the fields used commonly in benign traffic. These unknown HTTP-based malicious traffic is highly similar to benign traffic, usually can bypass the detection systems. Therefore, it is difficult but necessary to detect HTTP-based malicious communication behavior, especially unknown. This motivates researchers to pursue advanced detection techniques.

The key to detecting unknown HTTP-based malicious communication behavior is to improve the generalization ability of detection methods, i.e., the ability to discover unknown attacks by the known. There are two main challenges for detecting:

1) ***Feature extract.*** Traffic feature of HTTP-based malicious communication behavior is complex. However, many detection methods ([44, 33], *etc.*) rely on feature rules and expert knowledge so that they are difficult to grasp the essential laws of HTTP-based malicious communication behavior. The establishment of artificial feature engineering in a single experimental environment will cause the model to over-fit, which limits the generalization of method. This makes it difficult for a method to detect unknown HTTP-based malicious communication traffic. Namely, a method may perform well in assigned datasets, but poor in other large-scale datasets and real traffic environments.

2) ***Experimental dataset.*** The data scale associated with HTTP-based malicious communication behavior is relatively small. Many works ([1, 34], *etc.*) rely on private experimental data (the collection period is short and the collection points are limited), or rely on public malicious traffic datasets. However, due to the issue of timeliness, it is difficult for those datasets to cover all forms of such malicious traffic that occurs in the future.

In this paper, to cope with the above challenges in the detection of unknown HTTP-based malicious communication behavior, we propose an HTTP-based Malicious Communication traffic Detection Model (HMCD-Model). The main contributions are as follows:

- **We analyze the HTTP-based malicious communication behavior from the perspective of adversary.** We employ WGAN-GP [16] to synthesize Generated Adversarial Flows (GAFs) with maliciousness, compliance, covertness and multiformity for data enhancement. GAFs look highly similar to benign flows, which can be used as a supplement to labeled data and improve the generalization of HMCD-Model.
- **We propose a prototype system composed of a hybrid neural network to extract the hierarchical spatial-temporal features of HTTP traffic from packet level and flow level.** In this system, CNN is used to extract the spatial features of a packet, and LSTM is used to extract the temporal features of a flow. In addition, statistical features are also extracted hierarchically to improve detection performance.

And we discard miss-leading attributes (ip, url, *etc.*) that could easily lead to mis-judgment in data pre-processing.

- **We publish a dataset HMCT-2020 based on the real network environment**, which consists of large-scale HTTP-based malicious communication traffic and benign traffic during three years (2018-2020)<sup>1</sup>. There are about 76,760 malicious flows and 4,798,110 benign flows. HMCT-2020 can not only support our experiments but also help researchers further study HTTP-based malicious communication behavior.

Experimental results show that HMCD-Model has excellent detection performance. In HMCT-2020 dataset,  $F1$  is 99.46%(+0.19, -0.30), and  $FPR$  is 0.48%(+0.20, -0.34). For generalization (Taking the data in HMCT-2020(18) as the training set and the data in other datasets as the test set), our model has obvious advantages compared with the representative works in HMCT-2020 and the public dataset CIC-IDS-2017[35]. In addition, we also collect malicious traffic generated by malware and a large amount of benign background traffic from the real world, In the comparative experiment, the  $F1$  and  $FPR$  of HMCD-Model can reach 83.66%(+2.15, -3.79) and 2.57%(+3.26, -1.71), which are also better than baselines and other methods[44, 33]. The results prove that our method has a stronger ability to discover unknown malicious communication behavior.

The remainder of this paper is organized as follows. Section 2 introduces the related work. In Section 3, we analyze the HTTP-based malicious communication behavior. Feature analysis of HTTP traffic is in Section 4. Subsequently, Section 5 introduces the composition of HMCD-Model. In Section 6, we evaluate our method and show the relevant experimental results. Finally, we discuss and summarize in Section 7 and Section 8, respectively.

## 2. Related Work

In cyber security, malicious behavior detection is a hot issue[22]. We detect unknown HTTP-based malicious communication behavior based on deep learning, which belongs to the field of intrusion detection. Next, the related research status will be introduced.

### 2.1. Malicious Behavior Detection

Malicious behavior detection methods used in Intrusion Detection Systems (IDSs) [29] can be divided into feature detection and anomaly detection. Feature detection [3], also called misuse detection, fits the behavior patterns of known attacks and establishes a corresponding feature behavior database. Anomaly detection [5], also called behavior detection, mainly builds a feature database by fitting the characteristics of benign network behavior. Anomaly detection is slightly weaker than feature detection when detecting known attacks. However, anomaly detection can detect 0-day attacks more effectively. And this is very important for network security, because the network environment is becoming more and more complex, new 0-day attacks are constantly occurring, and a method that can effectively detect new attacks is necessary.

---

<sup>1</sup>The published dataset can be found at <https://github.com/BitBrave-Xie/HMCD-Model>.

### 2.1.1. Malicious Behavior Detection Based on Feature Selection and Machine Learning

Many works are based on some public network datasets (ISCX-2012 [37], KDD CUP 99 [40], CIC-IDS-2017, *etc.*), and then malicious behavior detection methods combining feature selection and machine learning are employed.

Aburomman *et al.* [1] review intrusion classification algorithms based on commonly used methods in the field of machine learning. In particular, considering integration methods of homogeneous and heterogeneous types, various integration and hybrid technologies are studied.

Wang *et al.* [44] propose an effective and automatic malware detection method using the text semantics of network traffic, treating each HTTP flow generated by software as a text document and processing it through N-gram to extract text-level features. Then, an automatic feature selection algorithm based on chi-square test is used to identify meaningful features, and these features are used to establish a support vector machine (SVM) classifier [42] for malicious behavior detection.

Salo *et al.* [33] propose an ensemble classifier based on SVM, Instance-based learning algorithms (IBK) [2], and multilayer perceptron (MLP) [46] for intrusion detection, which combines the approaches of Information Gain (IG) and Principal Component Analysis (PCA). The experimental results in datasets ISCX 2012, NSL-KDD [43] and Kyoto 2006+ [39] show that IG-PCA-Ensemble can learn more key features, and its classification accuracy, detection rate and false alarm rate are better than most of the existing advanced methods.

There are also many other studies about malicious behavior detection. Some methods are mainly based on collected experimental datasets. Wang *et al.* [45] propose BotMark for bot-nets detection based on flow-based and graph-based network traffic behavior. Du *et al.* [13] use SVM to differentiate the two types of anomaly in the mixed traffic. Shrestha *et al.* [38] use SVM for covert channel detection. Others methods are mainly based on public intrusion detection datasets. Zhou *et al.* [48] propose a heuristic dimension reduction algorithm CFS-BA-Ensemble based on feature selection and ensemble learning technology. Selvakumar *et al.* [34] deploy filters and wrappers based on the firefly algorithm in the feature selector, and use C4.5 and BN to classify in KDD CUP 99 dataset. Hajisalem *et al.* [18] propose a hybrid classification method based on Artificial Bee Colony (ABC) [23] and Artificial Fish Swarm (AFS) [17] algorithms. The method is superior to traditional machine learning methods in the NSL-KDD dataset and the UNSW-NB15 dataset [30].

### 2.1.2. Malicious Behavior Detection Based on Deep Learning

Deep learning is widely researched and applied in the field of intrusion detection due to its powerful feature extraction capabilities [25]. After the traffic is simply pre-processed, the neural network can automatically extract features and do not require researchers to put more effort into establishing feature engineering.

Chowdhury *et al.* [9] extract outputs from different layers in CNN and implement a linear SVM and 1-nearest neighbor classifier for few-shot intrusion detection. Kim *et al.* [24] propose a multi-mode deep learning method for malware detection. Caviglione *et al.*

[4] use neural networks and decision trees to detect malware using covert channels. Du *et al.* [12] use LSTM to perform malicious behavior detection.

## 2.2. Traffic Generation based on GAN

Currently, many researchers use deep learning technology to synthesize traffic data, which is used to bypass detection systems or enrich experimental data. Attempts to introduce the GAN to this field have shown promise. Zingo *et al.* [49] propose the "GAN vs Real (GvR) score", a task-based metric which quantifies how well a traffic GAN generator informs a classifier compared to the original data. Experiments show that it is possible to train accurate traffic anomaly detectors with GAN-generated network traffic data based on GvR.

Li *et al.* [26] propose a dynamic traffic camouflaging technique, coined FlowGAN, to dynamically morph traffic feature as another "normal" network flow to bypass Internet censorship. The core idea of FlowGAN is to automatically learn the features of the "normal" network flow, and dynamically morph the on-going traffic flows based on the learned features by GAN. Experimental results on a dataset involving 10,000 realworld flows show that the effectiveness and the efficiency of FlowGAN.

Ring *et al.* [31] generate flow data based on GANs and propose three different flow-based data pre-processing methods in order to convert them to continuous values. On this basis, a network traffic evaluation method based on domain knowledge definition quality test is proposed. Experiments on the CIDDS-001 dataset [32] show that two of the three methods can generate high-quality data.

Lin *et al.* [27] propose a framework based on GANs, namely IDSGAN, to generate adversarial samples to evade the detection of IDSs. IDSGAN uses generators to convert original malicious traffic into hostile malicious traffic and uses discriminators to simulate a black box detection system.

Cheng *et al.* [7] propose a GAN method for creating network traffic data at the ip packet layer. It prove feasibility in the generation of real traffic flows such as ICMP Pings, DNS queries, and HTTP web requests. Experiments show that the generated packets can be successfully transmitted through the Internet and the corresponding response can be obtained.

Jan *et al.* [21] propose request data synthesis method to synthesize unseen (or future) robot behavior distribution. The synthesis method has distributed perception capabilities and uses two different generators in GAN to generate data for clustering regions and outliers in the feature space.

Hao *et al.* [19] propose a GAN-based data augmentation method. The features of flow-based network traffic are first preprocessed to fit the GAN, and then the Earth-Mover (EM) distance is employed to capture the distribution of low-dimensional subspace data, while an encoder structure is added to learn latent space representations to enhance the vanilla GAN. They construct an imbalanced dataset based on a real-world dataset and compare it with other methods, obtaining better performance in terms of recall, F1 score and AUC.

Cheng *et al.* [8] propose Attack-GAN based on the structure of SeqGAN [47], to generate domain-constrained adversarial network traffic at the packet level. Specifically, adversarial packet generation is formulated as a sequential decision process. In this case, each byte in

the packet is considered a token in the sequence. The generator’s goal is to choose a token that maximizes its expected final reward. Generated network traffic and benign traffic are classified by black box IDS. The prediction results of IDS are fed into the discriminator to guide the update of the generator. Experimental results verify that the generated adversarial examples are able to deceive many existing black-box IDSs.

### 2.3. Analysis and Summary

For malicious behavior detection, there are two main challenges in the above methods for detecting unknown HTTP-based malicious communication behavior based on the discussions above methods. First, these methods have a valuable reference for feature modeling, but it is difficult to directly apply to the feature extracting of HTTP-based malicious communication behavior under adversarial conditions. For instance, some unstable features (ip, url, *etc.*) used before may be invalid. Moreover, hierarchical spatial-temporal features have not been considered too much. Second, these methods are usually tested on a single small-scale dataset, which cannot fully verify their generalization ability.

For traffic generation based on GAN, current methods for generating traffic samples are basically imitating specific datasets (CIDDS-001, *etc.*) or specific formats to generate various flow-based statistical data (number of packets, duration of the flow, *etc.*), not to generate real traffic that can be transmitted on the network, especially HTTP-based malicious communication behavior. This makes various traffic generation technologies have great limitations and can only be applied to a specific experimental scene.

Therefore, we build a hybrid neural network model for HTTP-based malicious communication behavior detection. Then, we propose a generation algorithm with good generality to synthesize GAFs. In addition, we construct a well-represented HTTP-based malicious communication traffic dataset to verify our model under different experimental conditions.

## 3. Analysis on HTTP-based Malicious Communication Behavior

We show the general process of HTTP-based malicious communication behavior, as shown in Fig. 1. An HTTP-based malware attack can be divided into four phases: implantation phase, incubation phase, communication phase and execution phase. 1) During the implantation phase, the adversary scan the victim device, or the victim accesses the StepStone/C&C server. Then, an HTTP-based malware script is covertly downloaded to the victim device. 2) During the incubation phase, usually, the malicious script sleeps for a period of time and enters the incubation phase to avoid being detected. 3) Then, the malware script enters the communication phase and sends on-line packets to contact the StepStone/C&C server. 4) Finally, the malware script enters the execution phase and begin to carry out various local actions according to instructions of the adversary.

HTTP-based malicious communication behavior can be effectively detected in the communication phase according to the instruction and the type of malicious interaction, because this is the phase where the characteristics of malicious communication behavior are most obvious. During the communication phase, defenders can analyze the traffic generated by

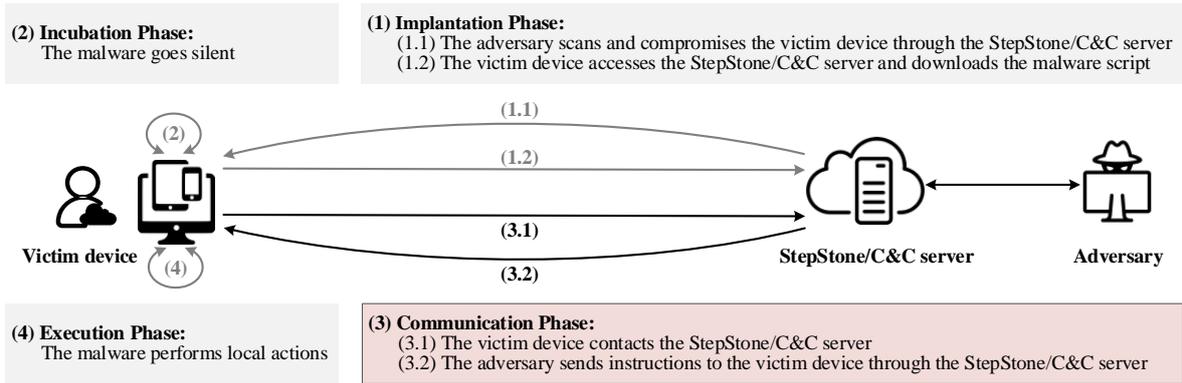


Figure 1: General process of HTTP-based malicious communication behavior.

malware, then, detect HTTP-based malicious communication behavior and find the adversary.

However, experienced adversaries usually construct HTTP-based malicious communication traffic by imitating the network behavior of benign users to evade the detection of defenders, making it difficult to distinguish from benign traffic. For instance, fig. 2 shows a flow consisting of two packets that belongs to HTTP-based malicious communication behavior. And we can divide the packet into three different components: 1) *Malicious components*, contain malicious content and has actual attack significance. For instance, the part enclosed by dotted line “jk?c=2&p=f4xZ24H4EPu\_hGnx7oUJh0cvXvXo50j9\_VnCX+zIl=&k1” in Fig. 2; 2) *Fixed components*, are the components with a fixed format that may not exist in an HTTP packet but cannot be changed; 3) *Covert-ness components*, are the other components except malicious and fixed component. Based on these three components, an adversary can synthesize a variety of malicious traffic to evade the detection of defenders by imitating benign traffic (such as embedding different instructions in the malicious components).

```

GET
jk?c=2&p=f4xZ24H4EPu_hGnx7oUJh0cvXvXo50j9_VnCX+zIl=&k1
HTTP/1.1
Accept: */*
User-Agent: MeDcore
Connection: Keep-Alive
Content-Type: text/xml; charset="UTF-16LE"
Host: [REDACTED]

HTTP/1.1 200 OK
Server: Tengine
Date: Tue, 25 Jul 2019 14:27:35 GMT
Content-Type: application/zip
Content-Length: 258
Connection: close
\x08\x00\x02\x00\x01\x0b\x02\x00\x00\x00\x07succeed\x0c\x00r\x02
\x00\x03\x00\x00 .....

```

Figure 2: A malicious flow consisting of two packets (a request and a response) generated by HTTP-based malicious communication behavior during the communication phase.

In this paper, we focus on the scene that an HTTP-based malware attack the victim devices and we need find a method to detect it during the communication phase. Therefore, we build a hybrid neural network to extract features of malicious traffic, and design

a traffic generation algorithm based on the four characteristics of HTTP-based malicious communication traffic.

#### 4. Feature Analysis of HTTP-based Malicious Communication Traffic

We define a flow as a sample. Flow is full-duplex in application layer, including request and response packets with the same quintuple ( $src\_ip$ ,  $src\_port$ ,  $dst\_ip$ ,  $dst\_port$ ,  $TCP$ ) over a period of time. A flow has multiple packets. Each packet is composed of payload and different fields, which has hierarchical structures. Therefore, we divide a flow into packet level (Pkt-level) and flow level (Flow-level) to extract features from different hierarchies.

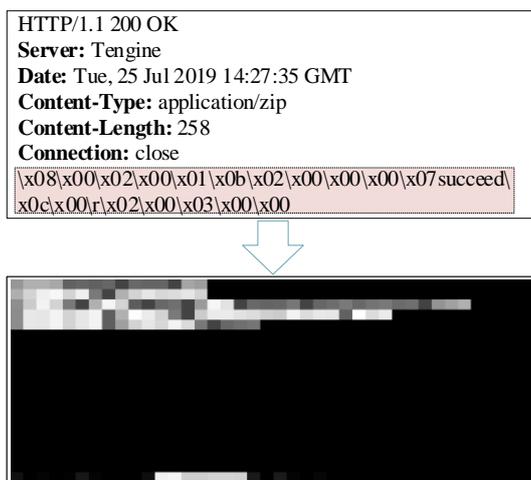


Figure 3: A example of converting a packet(top) of a flow into a two-dimensional image(bottom).

##### 4.1. Feature Analysis of Packet

**Text Feature:** The text content of an HTTP packet is composed of a payload and different header fields. The feature information of HTTP-based malicious communication behavior is usually hidden in these contents. For detecting unknown malicious communication behavior, we consider that the miss-leading contents (ip, url, *etc.*), will bring false positives. For instance, the host field of the packet can be changed to evade detection when the malicious communication information embedded in a packet is consistent. Therefore, we drop these miss-leading contents. Then, a packet is processed as a two-dimensional image with one channel, as shown in Fig. 3.

**Statistical Features:** Malicious and benign packets are different in many statistical values. For instance, an adversary usually uses fewer fields for simplicity. The length of the domain name used by the adversary is longer and unreadable due to the occupation of the domain name space. Therefore, we propose statistical feature engineering combined with the raw data, to further improve the performance of the detection model. The statistical features at Pkt-level are shown in Tab 1. RFC1998 [6] recommends that web services use 47 field lines for HTTP-based communication, but most web services do not use that much. Therefore, we consider counting the features of 18 fields. It will be discarded if exceeded and be filled with 0 if missed.

Table 1: Statistical features at Pkt-level

Type	Position
packet type	0
length of url or state description	1
protocol version	2
lines of fields	3
lengths of fields name	4-21
lengths of fields value	22-39
length of payload	40

#### 4.2. Feature Analysis of Flow

HTTP-based malicious communication traffic consists of multiple packets in a flow. The features provided by a single packet are limited, and analysis based on the entire flow can get more information. Generally, there are more packets in a malicious flow compared with benign traffic behavior. Usually, the bytes of response are larger but the bytes of request are smaller. Therefore, we perform feature analysis at Flow-level. Statistics such as the number of packets and length sequences are used as part of statistical feature engineering, as shown in Tab 2. The number of packets of a flow will not exceed 50, if it is exceeded, it will be discarded, or if it is missing, it will be filled with 0.

Table 2: Statistical features at Flow-level

Type	Position
count of request pkts	0
count of 'get'	1
count of 'post'	2
count of 'head'	3
count of 'options'	4
count of other requests	5
count of response pkts	6
count of '1XX'	7
count of '2XX'	8
count of '3XX'	9
count of '4XX'	10
count of '5XX' and others	11
count of other responses	12
mean of pkt bytes	13
seq of pkts bytes	14-63

## 5. HMCD-Model Methodology

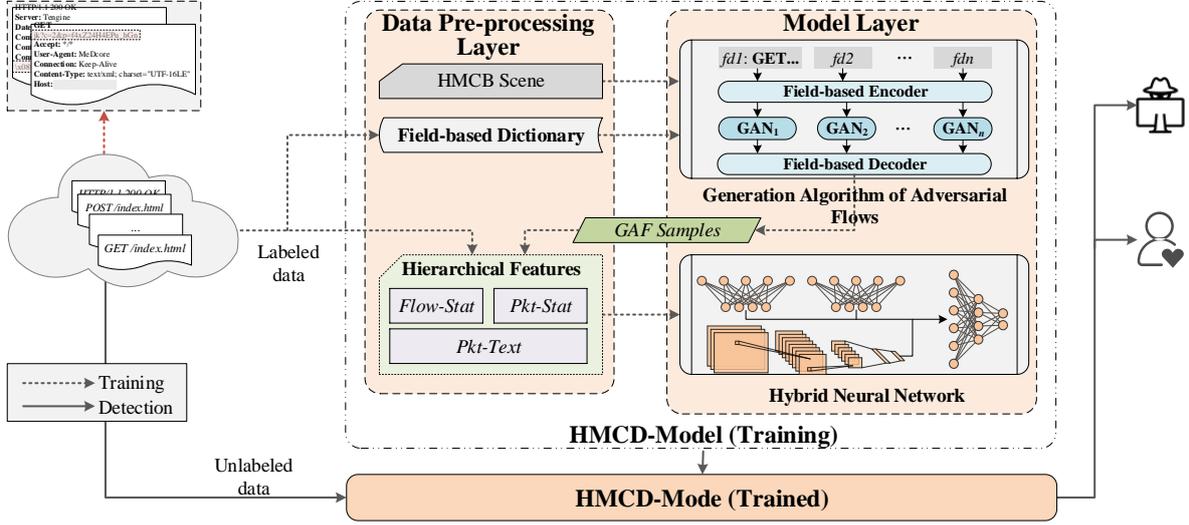


Figure 4: Training and detection process for HTTP-based malicious communication traffic (HMCSB Scene: HTTP-based malicious communication behavior scene).

### 5.1. Overview

As shown in Fig. 4, we build HMCD-Model. The training and detection process of the model is composed of two different stages.

**Generation of Adversarial Flows:** Generate HTTP-based malicious communication traffic for data enhancement. Based on a field-based dictionary, we use WGAN-GP to insert malicious content into benign traffic packets to generate malicious flows, that is Generated Adversarial Flows (GAFs), which are used to enrich the multiformality of training samples.

**Training and Detection:** Build a hybrid neural network. First, the labeled traffic data from the real world is used for model training, and GAFs are also added to the training set to improve the generalization of the model. Then, the trained HMCD-Model fixes the parameters, and performs unknown HTTP-based malicious communication behavior detection in various experimental scenarios.

Intuitively, we first introduce the hybrid neural network used to extract features, and then introduce the adversarial flows generation algorithm.

### 5.2. Hybrid Neural Network based on Hierarchical Spatial-temporal Features

The hybrid neural network based on hierarchical spatial-temporal features is shown in Fig. 5, and the basic unit of its detection is flow. In the spatial, the text features of the packet have structural information. So we use CNN to extract content and structural features. In the temporal, a flow usually has multiple packets, and there is natural timing between the packets. So we use LSTM to extract its sequence features.

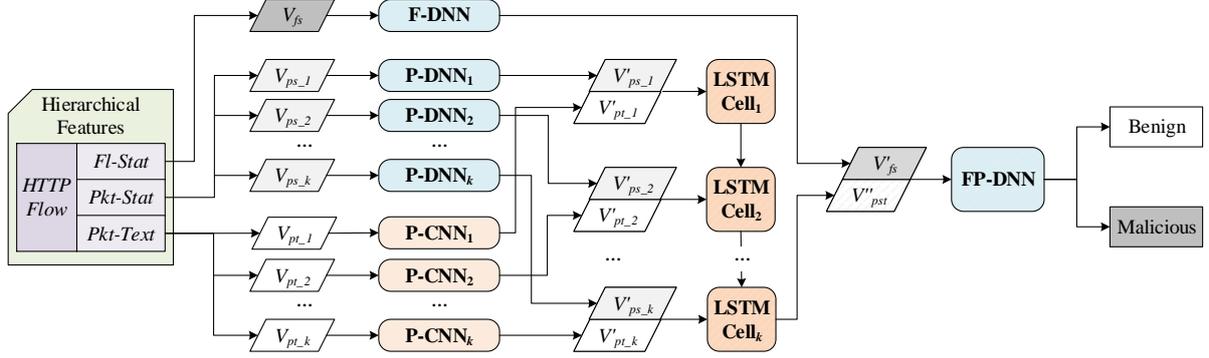


Figure 5: Structure of the hybrid neural network ( $V_{fs}$ : normalized Flow-Stat vector;  $V_{ps-*}$ : normalized Pkt-Stat vector;  $V_{pt-*}$ : normalized Pkt-Text vector;  $V'_*$ ,  $V''_*$ : intermediate variables;  $k$ : the packet size in an HTTP flow).

### 5.2.1. Feature Extract at Pkt-level

A packet has content and structure information. We cut a raw packet into a two-dimensional image with one channel to form text feature of the packet. In addition, we extract packet statistics to obtain more comprehensive packet information.

**Text Feature Extract of Packet:** The hybrid neural network uses CNN to process packet text information (Pkt-Text). CNN exploits convolution and pooling operations, which can perform translation-invariant classification of input information according to its hierarchical structure [14]. The single element output of the convolution layer corresponds to a matrix feature map input, as shown in Eq (1), where  $Z^l$  is input feature,  $w^{l+1}$  is convolution kernel,  $f_{conv}$  is the size of kernel and  $s_0$  is the stride. Convolution kernel parameters are shared, and the connection between layers are sparse, which allow it to effectively extract features with a small amount of calculation.

The pooling layer performs feature selection and information filtering in the input features. As shown in Eq (2), where  $f_{pool}$  is the size of pooling. In this paper, we use maximum pooling ( $p \rightarrow +\infty$ ) to retain the most significant features.

$$\begin{aligned}
 Z^{l+1}(i, j) &= [Z^l \otimes w^{l+1}](i, j) + b \\
 &= \sum_{k=1}^{K_l} \sum_{x, y=1}^{f_{conv}} [Z_k^l(s_0 i + x, s_0 j + y) w_k^{l+1}(x, y)] + b
 \end{aligned} \tag{1}$$

$$Z_k^l(i, j) = \left[ \sum_{x, y=1}^{f_{pool}} Z_k^l(i + x, j + y)^p \right]^{\frac{1}{p}} \tag{2}$$

**Statistical Feature Extract of Packet:** The statistics of a single packet is a 41-dimensional vector (Pkt-Stat), and there is no obvious characteristic relationship between the components. Therefore, we exploit DNN to process the Pkt-Stat and integrating them into lower latitude feature space.

### 5.2.2. Feature Extract at Flow-level

Packets of a flow have natural timing relationship according to the transmission timing. These packets form sequence features with a consistent format after packets are processed by CNN and DNN. Therefore, we extract sequence and statistical features in a flow at Flow-level.

**Sequence Feature Extract of Flow:** Packets of the same flow are serially transmitted on a timeline to form sequence data, which are suitable for processing using RNN. And the packets in the flow may have a long dependency due to delays and re-transmissions. LSTM can alleviate the long-term dependency problem [20]. Therefore, we use LSTM instead of traditional RNN to extract the sequence features of flow. It uses gate structure to selectively remember and forget information. Each gate is composed of an activation layer and a pointwise operation. By choosing to store information for subsequent processing, information can be transmitted further along the timing chain. Sundermeyer *et al.* [41] show that the most important component in LSTM is the forget gate. As shown in Eq (3), the forget gate,  $f_t$ , decides the proportion of the retained information according to the stored information  $h_{t-1}$  and the input  $x_t$  at the current stage. Followed is the input gate,  $i_t$ . And it decides how much new information to add, as shown in Eq (4).

$$f_t = \delta(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (3)$$

$$i_t = \delta(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (4)$$

**Statistical Feature Extract of Flow:** The statistics of a flow (Flow-Stat) is a 64-dimensional vector. There is no obvious relationship between the components. Therefore, we exploit DNN to process the vector. The input features are dimensional reduced and integrated into low-latitude feature vectors, which combined with LSTM output and then entered into subsequent neurons.

### 5.2.3. Feature Aggregation

FP-DNN takes the sequence and statistical features as input, which processed by LSTM and DNN. It exploits fully-connected layers for processing and outputs the final detection result. The final detection result of HMCD-Model is output after feature reduction and nonlinear transformation.

### 5.2.4. Time Complexity of Hybrid Neural Network

Time complexity is the number of executions of each operation in an algorithm or model, and is called operation frequency or time frequency. The time complexity of the neural network is the number of operations for each basic operation of the model in training and testing.

The time complexity,  $T_{hnn}$ , consists of the time of each local networks during training and testing, including feature extract and feature aggregation, as shown in Eq (5), where  $T_{dnn}$  is the sum of  $T_{p-dnn}$ ,  $T_{f-dnn}$  and  $T_{fp-dnn}$ ,  $N$  is the sample size. When model structure (the size of packet, *etc.*) are determined, its time complexity is linearly related to the sample size, that is,  $O(N)$ . By the way, in this hybrid neural network, the number of neurons in

each component is small to to improve the training and detection speed (convolution kernel, *etc.*). And some local networks can work simultaneously (*e.g.* text feature and statistical feature extract of packet, sequence feature and statistical feature extract of flow), and most neural networks are matrix operations, so there are many parallel computing strategies that can further accelerate the model.

$$\begin{aligned}
 T_{hnn} &= N \cdot (T_{cnn} + T_{lstm} + T_{dnn}) \\
 &\sim O(\beta \cdot N) \sim O(N)
 \end{aligned}
 \tag{5}$$

### 5.3. Generation Algorithm of Adversarial Flows

In an adversarial environment, the HTTP-based malicious communication traffic collected from Internet is difficult to cover all the specific forms. Therefore, we design a generation algorithm for generating adversarial flows. The algorithm first builds a dictionary of packet fields (field-based dictionary) based on the real traffic dataset, and then, generates corresponding packets, and combines the packets into flow.

In this paper, each of our generated adversarial flow consists of a request packet and a response packet. The process of generation algorithm is described in detail below.

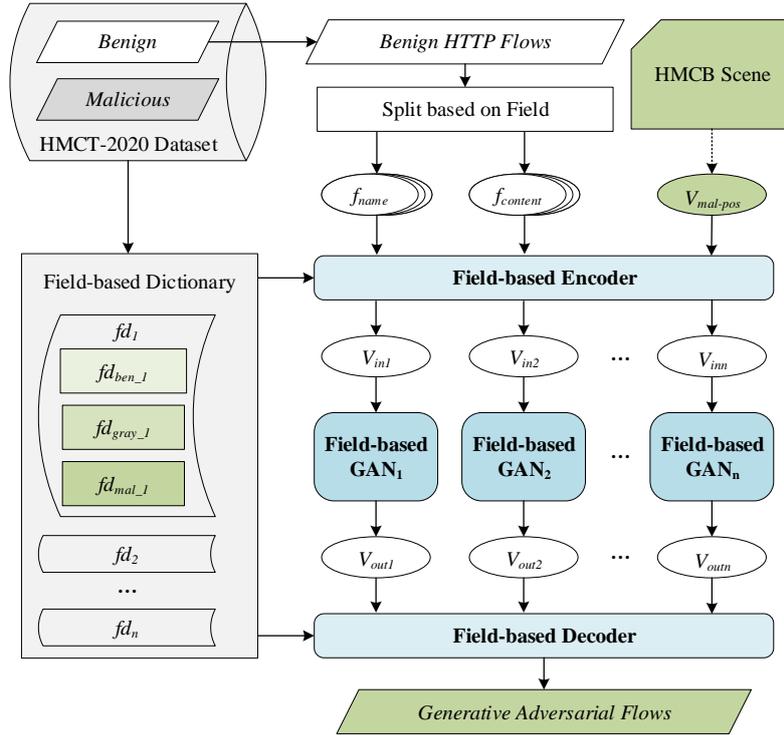


Figure 6: Process of generating HTTP-based malicious communication traffic (HMCB Scene: HTTP-based malicious communication behavior scene).

### 5.3.1. Generation Process of GAFs

The adversary always hides the malicious content in certain packet fields in the scene of HTTP-based malicious communication behavior. Therefore, we first generate each content in a flow based on the fields of a packet. Then, we splice these field-based contents to synthesize malicious packets, and then splice these packets to flows.

The process of generating HTTP-based malicious communication traffic is shown in Fig. 6. The process is divided into: 1) We analyze the traffic and generate corresponding benign and malicious dictionaries. 2) We randomly select the malicious content in the malicious dictionary and determine the location where it appears in HTTP benign flows, and finally obtain the corresponding field code  $V_{in}$ . 3) We input multiple  $V_{in}$  of a packet to the corresponding GAN for training, and output the generated field code  $V_{out}$ . 4) According to the HTTP specification, multiple  $V_{out}$  are decoded to packets and spliced into complete malicious communication flow.

**Field-based Dictionary:** We build a field-based dictionary,  $field\_dict$ , as shown in Eq (6).  $field\_dict$  is a two-level dictionary and contains three second-level dictionaries (malicious dictionary ( $fd_{mal}$ ), gray dictionary ( $fd_{gray}$ ) and benign dictionary ( $fd_{ben}$ )). HTTP packets have fixed structure, each row is independent and has specific meaning. Therefore, we divide packets by row, and each row is split into a fixed field-based name  $f_{name}$  (“GET”, “Accept”, *etc.*) and a remaining fixed field-based content  $f_{content}$  (may contain malicious information). The key of the first-level dictionary is  $f_{name}$  (“Date”, *etc.*) and the value is a second-level dictionary. According to the syntax and semantics of the  $f_{content}$ , we divide it into several words (“text”, “xml”, *etc.*) by special characters (“,”, “:”, *etc.*).

$$field\_dict = \left\{ \begin{array}{l} f_{name.1} : fd_{mal.1} + fd_{gray.1} + fd_{ben.1} \\ f_{name.2} : fd_{mal.2} + fd_{gray.2} + fd_{ben.2} \\ \dots \\ f_{name.n} : fd_{mal.n} + fd_{gray.n} + fd_{ben.n} \end{array} \right\} \quad (6)$$

For instance, we show an example in Fig. 2. These two packets are divided into 11 pairs of  $f_{names}$ , and  $f_{contents}$ . In the first one,  $f_{name}$  is “GET”, and  $f_{content}$  (“jk?c=2&p=f4Z24...”) can be split into words (“jk”, “c”, “2”, “p”, “f4Z24...”, *etc.*) by the characters like “?”, “=”, “&”. We use words that appear more than  $p$  times as keys to make the dictionary statistically stable and control the size of dictionary by adjusting the value of  $p$ . In addition, keys of  $fd_{mal}$  only appear in malicious HTTP packets, keys of  $fd_{ben}$  only appear in benign HTTP packets, and keys of  $fd_{gray}$  appear in both.

**Encoder:** According to the analysis about HTTP-based malicious communication behavior (see Section 3), we set the vector  $V_{mal-pos}$  to present the possible positions of packets in a flow where adversaries usually hide malicious content. As shown in Eq (7), where  $p_i$  is the replacement position,  $w_i$  is the malicious content word randomly selected from the keys in the  $fd_{mal}$ ,  $m$  is length of the  $V_{mal-pos}$ . Several parts of benign contents are replaced with malicious contents based on  $V_{mal-pos}$ . Then, the field-based contents through replacement are encoded as integer vectors ( $V_{in.1}$ ,  $V_{in.2}$ , *etc.*).

$$V_{mal-pos} = \{ \langle p_1, w_1 \rangle, \langle p_2, w_2 \rangle, \dots, \langle p_m, w_m \rangle \} \quad (7)$$

**Field-based GAN:** The Field-based GAN is responsible to generate malicious content which is highly similar to benign content. As shown in Fig. 7, Field-based GAN consists of two modules: the Generator ( $G$ ) and the Discriminator ( $D$ ). The structure of the  $G$  is as same as  $D$ , but their order is reversed.  $G$  is to imitate real field-based benign contents and generate fake, and  $D$  is to judge whether the generated contents by  $G$  are similar to the real benign contents.

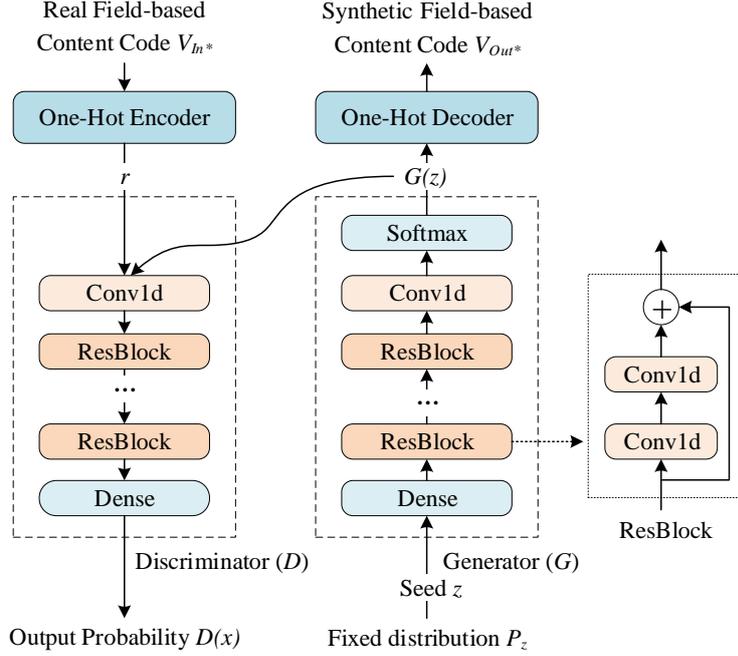


Figure 7: Structure of Field-based GAN.

In Field-based GAN, we introduce the idea of WGAN-GP [16], which makes Field-based GAN to generate diverse data. WGAN-GP defines a well-defined training procedure, which makes the generator easier to train than traditional GAN. It describes the distance between  $p_g$  and  $p_r$  by the Wasserstein distance. WGAN-GP can be considered a min-max game, as shown in Eq (8), where  $z$  is the input of  $G$  which is a random 1-D tensor  $z$  filled,  $R$  is the regular term and  $\mathbb{E}$  is the cross-entropy function. Through Eq (8), Field-based GAN continuously reverses iterative neural network training  $G$  and  $D$ , takes the random seed  $z$  as the input of  $G$ , output the content sample  $G(z)$ , and takes the  $V_{in}$  and  $G(z)$  (one-hot form) as the input of  $D$ . When all parameters converge, we decode the generated  $G(z)$  into the generated content  $V_{out}$ .

$$\min_G \max_D L(D, G) = \mathbb{E}_{z \sim P_z} [D(G(z))] - \mathbb{E}_{x \sim P_r} [D(x)] + R \quad (8)$$

**Algorithm Description:** The generation algorithm of adversarial flows based on the above-described generation process is shown in Algorithm 1. The input of the algorithm is the  $f_{name}$  and  $f_{content}$  of HTTP benign flows and the  $V_{mal-pos}$ . The output is GAFs. Algorithm guarantees *maliciousness* of GAFs through  $V_{mal-pos}$ . The strategy of dividing

packets based on fields ensures *compliance* of the generated traffic. Field-based GAN will imitate the characteristics of benign traffic when generating traffic, which ensures *covertiness* of the generated traffic. In addition, we use the WGAN-GP to optimize the training process and enhance *multiformity* of GAFs.

---

**Algorithm 1** generation algorithm of adversarial flows

---

**Input:** *benign-flows* : HTTP-based benign traffic.  $V_{mal-pos}$ : the vector that may contain malicious information.  $N$  : the number of samples that need to generate.

**Output:** *g-flows*: generated HTTP-based malicious communication flows.

**Step 1:** Dictionaries

Segment words and build dictionaries *field\_dict* from HMCD-2020 datasets

**Step 2:** Encoder

**for**  $f_{field}$  in *field\_dict* **do**

The content-string is separated from the corresponding fields of benign traffic

Encode the content-string as  $V_{code}$  based on *field\_dict*

**end for**

**Step 3:** Field-based GAN

**for**  $t=0$  in number of iterations **do**

The *seed*( $z$ ) is transformed into the sample  $G(z)$  with the same dimension as  $V_{code}$

Convert  $G(z)$  to  $V_{mal-code}$  based on  $V_{mal-pos}$

Take  $V_{code}$  and  $V_{mal-code}$  as inputs to  $D$

Update  $D$  model and  $G$  model by Eq (8)

**end for**

**Step 4:** Generated Adversarial Flows (GAFs)

Decode the generated  $G(z)$  and merge it into a malicious communication flow *g\_flow*

*g-flows* =  $N$  generated HTTP-based malicious communication flows

---

### 5.3.2. Time Complexity of Generation Algorithm

The algorithm traverses the input data once in building the dictionary and encoding stage. In the generation stage, the time complexity ( $T_{ga}$ ) is shown in Eq (9), where  $N$  is the sample size,  $F$  is the number of Field-based GAN,  $C_l$  is the number of convolution kernels of  $l$ -th layer,  $M$  is the size of the feature map, and  $K$  is kernel size,  $D_{in}$  and  $D_{out}$  are the input and output of the fully-connected layer. All the parameters are constants after the model structure is determined. Therefore, the time complexity of the entire generation algorithm ( $T_{ga}$ ) can be regarded as  $O(N)$ .

$$\begin{aligned}
 T_{ga} &\sim O\left(N \cdot F \cdot \sum_{l=1}^L (M_l^2 \cdot K_l^2 \cdot C_{l-1} \cdot C_l + 2 \cdot D_{in} \cdot D_{out})\right) \\
 &\sim O(\alpha \cdot N) \sim O(N)
 \end{aligned} \tag{9}$$

## 6. Experimental Evaluation

### 6.1. Datasets

#### 6.1.1. HMCT-2020 Dataset

Currently, there are no large-scale public intrusion detection datasets specifically for HTTP-based malicious communication behavior detection, the relevant datasets that contain HTTP-based malicious communication traffic have limited specific attack forms. Therefore, we publish the dataset HMCT-2020 and the hashing technology is used for data masking to protect privacy. HMCT-2020 can be divided into two datasets from the timeline, HMCT-2020(18) and HMCT-2020(19-20). The traffic in HMCT-2020(18) is mainly captured from the Internet during November 2018. The traffic in HMCT-2020(19-20) is captured from the Internet between July 2019 and March 2020. Finally, after data cleaning and application layer packet extraction, the captured traffic forms the dataset HMCD-2020, the details of which are shown in Tab 3. The dataset can be found in there<sup>2</sup>.

Table 3: Statistics on packet size and flow size in HMCT-2020

	Packet (in bytes)			Flow (in packets)				
		Malicious	Benign	Total		Malicious	Benign	Total
HMCT-2020 (18)	Count (in packets)	1.49M	11.48M	12.97M	Count (in flows)	35.58K	3.92M	3.95M
	Size	429.96M	6.5G	6.93G	Size	1.49M	11.48M	12.97M
	Min	15	12	12	Min	2	1	1
	Max	46.55K	10.23K	46.55K	Max	50	50	50
	Mean	289.25	566.33	534.58	Mean	41.77	2.93	3.28
HMCT-2020 (19-20)	Count (in packets)	68.77K	2.86M	2.93M	Count (in flows)	41.18K	882.0K	923.17K
	Size	29.6M	1.73G	1.76G	Size	68.77K	2.86M	2.93M
	Min	15	19	15	Min	1	1	1
	Max	6.32K	7.08K	7.08K	Max	50	50	50
	Mean	430.41	604.43	600.34	Mean	1.67	3.24	3.17

There are two sources of malicious traffic in HMCT-2020. One is to directly extract relevant malicious traffic from the network operator’s existing IDSs. The other is to actively capture traffic by running target malware in a sandbox. In HMCT-2020(18), a total of 35,583 flows are obtained, and in HMCT-2020(19-20), a total of 41,177 flows are obtained.

Benign traffic comes from the network gateway of our network security lab. After authorization, we deploy a traffic collection device at the gateway to collect traffic while ensuring

<sup>2</sup><https://github.com/BitBrave-Xie/HMCD-Model>

security and data privacy. Finally, HMCT-2020(18) contains the benign flows of 3,981,567, and HMCT-2020(19-20) contains the benign flows of 923,172.

We employ multiple ways to ensure the validity of traffic data labels in HMCT-2020. First, malicious traffic and benign traffic originate from different channels, and are naturally distinguishable so that they can be self-labeled. For instance, we class malicious traffic captured by partner operators’ existing IDSs and traffic generated by malware running in sandboxes as malicious, and traffic captured from lab network gateways as benign. Second, after cleaning the traffic data, we randomly sample from the preliminary labeled data to make artificial judgments to further increase the accuracy of data labeling. Finally, we construct the labeled dataset HMCT-2020.

### 6.1.2. Other Datasets

In this paper, the malicious traffic in the published intrusion detection dataset and the malicious traffic generated by the malware are collected to further test the generalization of model.

**CIC-IDS-2017:** CIC-IDS-2017[35] is an intrusion detection dataset with complete traffic published by the Canadian Institute for Cybersecurity at University Of New Brunswick, in 2018. It contains benign traffic and latest common attacks (Web Attack, Infiltration, *etc.*). We select all HTTP-based malicious traffic as positive samples and all the background traffic as the unknown benign traffic. The specific dataset details are shown in Tab 4.

Table 4: Statistics on packet size and flow size in CIC-IDS-2017

Packet (in bytes)				Flow (in packets)			
	Malicious	Benign	Total		Malicious	Benign	Total
Count (in packets)	131.1K	438.58K	569.68K	Count (in flows)	27.47K	106.23K	133.7K
Size	285.21M	580.78M	865.99M	Size	131.1K	438.58K	569.68K
Min	16	16	16	Min	2	1	1
Max	20.17K	24.71K	24.71K	Max	50	50	50
Mean	2.18K	1.32K	1.52K	Mean	4.77	4.13	4.26

**82-Malware-Traffic:** We collect HTTP traffic generated by 82 malware from an influential malware web<sup>3</sup>. The selection condition is malware with behavioral characteristics that meet the conditions of HTTP-based malicious communication behavior, such as *Gootkit*. In addition, we collect a large amount of background traffic from different Internet gateways as unknown benign traffic. The details about traffic can be found in our publish link (in Introduction). The specific dataset details are shown in Tab 5.

<sup>3</sup><http://www.malware-traffic-analysis.net/index.html>

Table 5: Statistics on packet size and flow size in 82-Malware-Traffic

Packet (in bytes)		Flow (in packets)	
	Malicious		Malicious
Count (in packets)	70.19K	Count (in flows)	3.19K
Size	22.94M	Size	70.19K
Min	15	Min	1
Max	7.41K	Max	50
Mean	326.81	Mean	21.98

## 6.2. Evaluation Metrics and Environmental Configuration

### 6.2.1. Evaluation Metrics

There are 4 basic metrics in the experiment. True-Positive ( $TP$ ), is the number of malicious samples classified as malicious. False-Positive ( $FP$ ), is the number of benign samples classified as malicious. True-Negative ( $TN$ ), is the number of benign samples classified as benign. False-Negative ( $FN$ ), is the number of malicious samples classified as benign.

Based on the above metrics, We use precision ( $P$ ), recall ( $R$ ) and false positive rate ( $FPR$ ) to evaluate a model’s performance in detecting HTTP-based malicious communication behavior, as shown in Eq(10) to Eq(12). In addition,  $F1$  is also used to verify the overall performance of a model in detecting malicious behavior. As shown in Eq(13),  $P_k$  and  $R_k$  are the precision and recall of the model in the k-th experiment. Because the same experiment will be repeated  $N_r$  times, it will produce multiple indicators of the same type, here we calculate the macro average to get the final result.

$$P = \frac{1}{N_r} \sum_{k=1}^{N_r} \frac{TP_k}{TP_k + FP_k} \quad (10)$$

$$R = \frac{1}{N_r} \sum_{k=1}^{N_r} \frac{TP_k}{TP_k + FN_k} \quad (11)$$

$$FPR = \frac{1}{N_r} \sum_{k=1}^{N_r} \frac{FP_k}{FP_k + TN_k} \quad (12)$$

$$F1 = \frac{1}{N_r} \sum_{k=1}^{N_r} \frac{2 \times P_k \times R_k}{P_k + R_k} \quad (13)$$

### 6.2.2. Environmental Configuration

The system environment is Ubuntu16.04 LTS. The hardware facilities are 16-core CPU and 128G memory. TensorFlow2.0 in Python3.7 is used to implement the model. To accelerate training and detection, 3 NVIDIA TITAN XPs are deployed on the server.

As a rule of thumb, we set some hyper-parameters for the model and related experiments. Field-based GAN has 7 convolutional layers and 2 full-connected layers. CNN has one convolutional layer with two  $2 \times 8$  neurons. LSTM cell uses 16 neurons. The structure of DNN is 10, 8, 2. The other parameters about HMCD-Model are shown in Tab 6.

Table 6: Parameter configuration of HMCD-Model during training

Type	Value
Loss function	Cross-Entropy
Activation function	ReLU
Optimizer	Adam
learning rate	1e-3
Batch size	128
Epochs	50

In addition, we set a sample consist the first two packets in a flow, and the size of a packet is  $20 \times 40$ . The number of repeated experiments is  $N_r = 5$ . During the training of model, the 5-fold cross validation is used.

For the experimental data, we conduct experiments by randomly sampling part of the traffic data from the datasets. Specifically, for the training set data, we randomly sample it in HMCT-2020(18), and for the test set data, we randomly sample it from different datasets according to the different requirements of the experiment. Tab 7 shows the details of the data composition of the training and test sets for each experiment. For instance, when detecting 82-Malware-Traffic, we select all malicious traffic. In addition, in order to ensure that the experimental results can truly reflect the performance of the model in the entire dataset, all data in each experiment are re-sampled randomly from the corresponding dataset.

For the number of GAFs in the training set, theoretically, the more GAFs, the higher the generalization of the HMCD-Model, but our preliminary experimental results show that it will reduce the performance of the model when the number of GAFs in the training set is too large. This is because a generated adversarial flow sample we generate consists of a request and a response, which can only imitate part of the HTTP-based real malicious communication behavior. The generated data can only change in a certain local feature space, so the excessive proportion of GAFs in the training set will lead to the performance of the model overfitting and reduce its generalization performance. In subsequent experiments, in the combination of training and test sets ep1 to ep4, we add 10,000 GAFs to the training set. That is, there are 30,000 malicious samples in the training set, consisting of 20,000 real traffic samples and 10,000 GAFs.

Table 7: Data composition of training set and test set in different experiments

	Training set			Test set		
	Data source	Malicious samples	Benign samples	Data source	Malicious samples	Benign samples
ep1	HMCT-2020(18)	20,000	50,000	HMCT-2020(18)	10,000	10,000
ep2	HMCT-2020(18)	20,000	50,000	HMCT-2020(19-20)	30,000	30,000
ep3	HMCT-2020(18)	20,000	50,000	CIC-IDS-2017	27,474	30,000
ep4	HMCT-2020(18)	20,000	50,000	82-Malware-Traffic	3,138	4,000

### 6.3. Ablation Study on Key Factors in HMCD-Model

We add statistical features to improve the detection performance of HMCD-Model and GAFs to improve the data multiformity. We verify these key factors in HMCT-2020 and observe corresponding improvements.

Table 8: Performance improvements of HMCD-Model due to the addition of statistical features and GAFs (ep2 in Tab 7)

	HMCT-2020(19-20) (ep2)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
HMCD-Model without Statistical Features	97.33 <sup>+0.75</sup> -0.62	99.37 <sup>+0.27</sup> -0.16	98.33 <sup>+0.34</sup> -0.32	2.73 <sup>+0.65</sup> -0.79
HMCD-Model without GAFs	94.63 <sup>+3.86</sup> -4.50	95.74 <sup>+2.68</sup> -5.96	95.15 <sup>+3.31</sup> -3.62	5.50 <sup>+5.15</sup> -3.99
HMCD-Model	<b>97.94</b> <sup>+0.7</sup> -1.51	<b>99.39</b> <sup>+0.18</sup> -0.25	<b>98.66</b> <sup>+0.38</sup> -0.89	<b>2.10</b> <sup>+1.57</sup> -0.73

#### 6.3.1. Statistical Features

We extract statistical features from Pkt-level and Flow-level. The addition of statistical features can make HMCD-Model get more comprehensive information and easier to extract the essential features of the flow. Traffic in HMCT-2020(18) is used training set, the experimental results are shown in Tab 8. All the indicators of HMCD-Model have been improved after adding statistical features.

#### 6.3.2. Generated Adversarial Flows

GAFs are used for data enhancement, to enrich the data multiformity of the training set. The experimental results are shown in Tab 8. After adding GAFs, the model can be increased to 98.66% in  $F1$ , and the  $FPR$  can be reduced to 2.10%. GAFs expand the

training set so that the model can learn more malicious behavior scenes, reduce over-fitting, and improve its generalization ability. We show a generated adversarial flow in Fig. 8, which consists of a request and a response.

```

POST /ctrl/agnostic/shards/14/&/v1/ctrl/v1/phpunit/as... HTTP/1.1
Host: steelcti1918035169cobweb.254.steelcn.cn152
Content-Type: jsonwwwUTFtext/dataformformtextcharset=wwwhtmlurle...
Content-Length: 430

HTTP/1.1 200 OK
Date: Wed2107, Apr032029531655
Server: 2 ./nginxCentOS)CentOSnginx
Content-Length: 430
Content-Type: jsonwwwUTFtext/dataformformtextcharset=wwwhtmlurle...

```

Figure 8: A generated adversarial flow composed of two packets (a request and a response).

In addition, we visualize the partial sample features distribution from the training set of HMCT-2020(18) and the corresponding GAFs feature distribution. The feature distribution based on t-SNE dimension reduction is shown in Fig. 9. HTTP-based malicious communication traffic and benign traffic are distributed in different feature spaces. Benign traffic can be grouped into a large cluster as a whole because they usually have some common characteristics. HTTP-based malicious communication traffic is distributed in multiple small clusters and the whole is relatively scattered. Because the behavior and purpose of different attackers represented by these flows are different, they often construct flows with extremely different characteristics to bypass different detection systems. The distribution of GAFs also shows the same phenomenon of small cluster distribution as HTTP-based malicious communication traffic. They are distributed in the feature space outside the benign traffic. This shows that our generation algorithm can effectively generate adversarial samples and fill the feature space of HTTP-based malicious communication traffic.

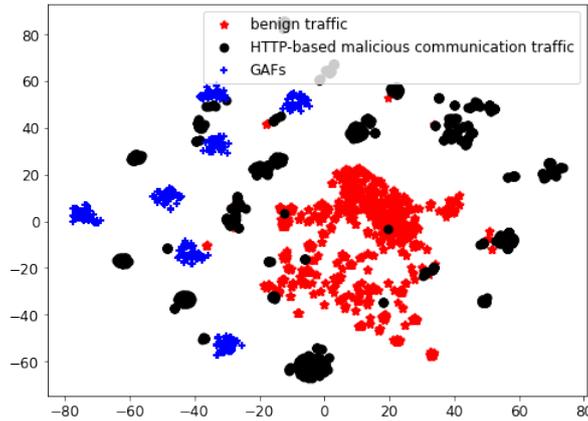


Figure 9: The feature distribution of HTTP-based malicious communication traffic, benign traffic and GAFs based on t-SNE dimension reduction.

#### 6.4. Generalization Comparison

We investigate research work related to HTTP-based malicious communication behavior and selected two well-represented works. Wang *et al.* [44] regard traffic as language and use

N-gram and SVM for detection, which can obtain 99.15% accuracy when detecting malware flows and can detect 54.81% for unknown malicious applications. Salo *et al.* [33] propose a novel hybrid dimensionality reduction technique for intrusion detection combining IG and PCA with an ensemble classifier based on SVM, IBK, and MLP, named IG-PCA-Ensemble. It can achieve 99+% performance in dataset ISCX-2012, NSL-KDD and Kyoto2006+.

We implement these two methods and select hyper-parameters of the two methods according to the trade-off strategy between detection performance and detection time in the experiments.

N-gram+SVM: Wang *et al.* We use two packets as the input feature for each flow. For other hyper-parameters, we follow the selection of the original paper under the reasonable time cost. We set the word length  $N = 1$  in N-gram and the feature number  $K = 600$  for each sample according to Wang’s experiment results.

IG-PCA-Ensemble: We set the number of neighbors of IBK to six, SVM is a linear kernel, and MLP has 128 neurons. It reads the first two packets for each flow, selects the top 17 most important information through IG, and reduces to 12-dimensional features based on PCA.

We compare the generalization of HMCD-Model, N-gram+SVM, and IG-PCA-Ensemble in a variety of experimental environments. In addition, we also compare HMCD-Model with classic baseline machine learning methods in 82-Malware-Traffic (see Appendix A).

#### 6.4.1. Generalization Comparison in Dataset HMCT-2020

We select data from the HMCT-2020(18) to form the training set and the test set. The results of experiments are shown in Tab 9. It shows that HMCD-Model has the best detection performance overall. In HMCT-2020(18), compared with N-gram+SVM, HMCD-Model has the same level of detection performance, but its  $F1 \approx 99.46\%$  is 13.8% higher than IG-PCA-Ensemble (the  $F1 \approx 85.66\%$ ).

Table 9: Performance comparison in dataset HMCT-2020(18) (ep1 in Tab 7)

	HMCT-2020(18) (ep1)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
N-gram+SVM	<b>99.79</b> <sup>+0.11</sup> <sub>-0.07</sub>	98.92 <sup>+0.62</sup> <sub>-2.18</sub>	99.35 <sup>+0.32</sup> <sub>-1.10</sub>	<b>0.21</b> <sup>+0.07</sup> <sub>-0.11</sub>
IG-PCA-Ensemble	83.92 <sup>+4.89</sup> <sub>-5.73</sub>	87.71 <sup>+1.93</sup> <sub>-2.71</sub>	85.66 <sup>+1.31</sup> <sub>-2.13</sub>	17.21 <sup>+7.79</sup> <sub>-6.50</sub>
HMCD-Model	99.52 <sup>+0.34</sup> <sub>-0.20</sub>	<b>99.39</b> <sup>+0.43</sup> <sub>-0.41</sub>	<b>99.46</b> <sup>+0.19</sup> <sub>-0.30</sub>	0.48 <sup>+0.20</sup> <sub>-0.34</sub>

In addition, we use the data of HMCT-2020(18) for training and the data of HMCT-2020(19-20) with more malicious types for testing to examine the generalization of the model against concept migration. The experimental results of the three methods are shown

in Tab 10. Obviously, HMCD-Model has significantly better detection performance than the other two methods. For instance, HMCD-Model’s  $FPR$  is only 2.1%, and the  $F1$  has a 35.67% and 26.07% improvement over N-gram+SVM, IG-PCA-Ensemble, respectively.

Table 10: Generalization comparison in dataset HMCT-2020(19-20) (ep2 in Tab 7)

	HMCT-2020(19-20) (ep2)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
N-gram+SVM	81.79 <sup>+1.83</sup> -1.07	51.24 <sup>+1.31</sup> -1.24	62.99 <sup>+0.73</sup> -0.84	11.42 <sup>+1.08</sup> -1.34
IG-PCA-Ensemble	75.92 <sup>+6.05</sup> -6.21	70.07 <sup>+4.68</sup> -4.87	72.59 <sup>+0.60</sup> -0.45	23.02 <sup>+9.46</sup> -8.68
HMCD-Model	<b>97.94</b> <sup>+0.70</sup> -1.51	<b>99.39</b> <sup>+0.18</sup> -0.25	<b>98.66</b> <sup>+0.38</sup> -0.89	<b>2.10</b> <sup>+1.57</sup> -0.73

The features of malicious traffic in HMCT-2020(19-20) have no major change in the overall spatio-temporal behavior compared to HMCT-2020(18) and only the attack information in packets is different. N-gram+SVM relies more on the content information in packets and works well when detecting malicious traffic similar to the training data in contents. When the information changes, its generalization performance will decrease significantly. IG-PCA-Ensemble relies more on the behavior characteristics of malicious traffic at Flow-level and is not sensitive to attack information, which results in its overall detection performance being centered. Compared with these two methods, HMCD-Model extracts hierarchical spatio-temporal of traffic features from Pkt-level and Flow-level respectively. Therefore, it shows excellent detection performance to resist concept migration.

#### 6.4.2. Generalization Comparison in Dataset CIC-IDS-2017

We selected the data in the dataset HMCT-2020(18) and CIC-IDS-2017 for generalization detection, one for training and the other for testing. Since the malicious communication traffic type in CIC-IDS-2017 is relatively single and the correspondingly constructed field-based dictionary cannot cover more malicious scenes, making it difficult for HMCD-Model to generate GAFs. Therefore, we use data in HMCT-2020(18) as the training set and data in CIC-IDS-2017 as the test set. The experimental results are shown in Tab 11. The precision of N-gram+SVM is  $P = 98.38\%$  and the  $FPR$  is 0.35, but other metrics are very low. This is because N-gram+SVM increases the threshold for determining a sample as malicious, which means a lot of false negatives. For instance, the  $F1$  of N-gram+SVM is 56.34%, and is 34.35% lower than HMCD-Model. In addition, N-gram+SVM and IG-PCA-Ensemble have large fluctuations in CIC-IDS-2017, which makes them easy to lose detection performance. Therefore, in general, our method has the best overall performance.

The detection results of these three methods have large fluctuations. The fluctuation of N-gram+SVM comes from the change of training data. N-gram+SVM is sensitive to content

Table 11: Generalization comparison in dataset CIC-IDS-2017 (ep3 in Tab 7)

	CIC-IDS-2017 (ep3)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
N-gram+SVM	<b>98.38</b> <sup>+1.31</sup> <sub>-2.44</sub>	51.28 <sup>+43.55</sup> <sub>-48.92</sub>	56.34 <sup>+40.86</sup> <sub>-51.73</sub>	<b>0.35</b> <sup>+0.53</sup> <sub>-0.30</sub>
IG-PCA-Ensemble	73.71 <sup>+15.68</sup> <sub>-62.25</sub>	75.04 <sup>+18.53</sup> <sub>-74.06</sub>	73.45 <sup>+17.98</sup> <sub>-71.65</sub>	10.50 <sup>+0.96</sup> <sub>-2.94</sub>
HMCD-Model	86.93 <sup>+12.03</sup> <sub>-11.13</sub>	<b>95.32</b> <sup>+3.14</sup> <sub>-2.80</sub>	<b>90.69</b> <sup>+6.77</sup> <sub>-5.03</sub>	15.37 <sup>+16.06</sup> <sub>-14.37</sub>

information of packets. It is difficult to extract features effectively when the difference between the content information of the training data and the test data exceeds a threshold. The fluctuation of IG-PCA-Ensemble mainly comes from the change of training data and the uncertainty of its own model. The fluctuation of HMCD-Model mainly comes from the initialization and update of parameters of the model itself. But it can be stabilized through multiple iterations of training, so the HMCD-Model is relatively more stable.

#### 6.4.3. Generalization Comparison in 82-Malware-Traffic

We capture relevant malicious traffic data from the Internet to verify the generalization performance of the model more comprehensively. We use data in HMCT-2020(18) as the training set and data in 82-Malware-Traffic as the test set. The experimental results are shown in Tab 12. Similarly, the precision of N-gram+SVM is  $P = 99.71\%$  and the  $FPR$  is 0.14, but other metrics are very low. This is because N-gram+SVM increases the threshold for determining a sample as malicious, which means a lot of false negatives. For instance, the  $F1$  of N-gram+SVM is 66.09%, and is 17.57% lower than HMCD-Model. In general, the HMCD-Model also has the best comprehensive performance and is more stable, with  $F1$  of 83.66% and  $FPR$  of 2.57%.

Table 12: Generalization comparison in 82-Malware-Traffic (ep4 in Tab 7)

	82-Malware-Traffic (ep4)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
N-gram+SVM	<b>99.71</b> <sup>+0.16</sup> <sub>-0.26</sub>	49.43 <sup>+0.29</sup> <sub>-1.09</sub>	66.09 <sup>+0.26</sup> <sub>-0.94</sub>	<b>0.14</b> <sup>+0.14</sup> <sub>-0.08</sub>
IG-PCA-Ensemble	94.18 <sup>+0.45</sup> <sub>-0.64</sub>	72.55 <sup>+14.03</sup> <sub>-20.92</sub>	80.81 <sup>+9.15</sup> <sub>-14.00</sub>	4.54 <sup>+1.44</sup> <sub>-1.61</sub>
HMCD-Model	96.62 <sup>+2.26</sup> <sub>-4.27</sub>	<b>73.77</b> <sup>+2.01</sup> <sub>-3.41</sub>	<b>83.66</b> <sup>+2.15</sup> <sub>-3.79</sub>	2.57 <sup>+3.26</sup> <sub>-1.71</sub>

The detection results of the three detection models in the real traffic environment has deteriorated compared with in HMCT-2020 and CIC-IDS-2017. This is because the real network environment is more complicated and the form of unknown malicious traffic is more variable. The traffic generated by malware is not all traffic related to malicious communication behavior, which leads to false positives.

#### 6.4.4. Comparison in Time Cost

We compared the detection time of HMCD-Model with N-gram+SVM and IG-PCA-Ensemble. Taking 10,000 samples as a unit, the detection time of each model is shown in Fig 10. IG-PCA-Ensemble spends the least time in detecting because its feature space is low-dimensional. However, there is no free lunch in the world. Although the detection time of IG-PCA-Ensemble is fast, its generalization performance is not satisfactory.

In addition, N-gram+SVM spends the second least time in detecting, because its feature space is relatively small and it uses linear SVM for detection. However, N-gram+SVM uses bag-of-words and Chi-square algorithms in pre-processing, which means it is difficult to arbitrarily increase the packet extraction content and feature space. Its storage space, pre-processing time and SVM convergence time will increase exponentially when there is too much input information.

The detection time of HMCD-Model is of the same magnitude as N-gram+SVM. In fact, the internal structure of HMCD-Model supports parallel processing, such as the processing of packet text and statistical features can be performed in parallel. The matrix operation characteristics of neural networks also enable HMCD-Model to make full use of the computing advantages of GPUs. More importantly, HMCD-Model maintains a good generalization detection ability, the detection time and the size of the feature space of the data are linear. Therefore, HMCD-Model is undoubtedly a more competitive and better comprehensive performance detection method from the current Internet big data environment.

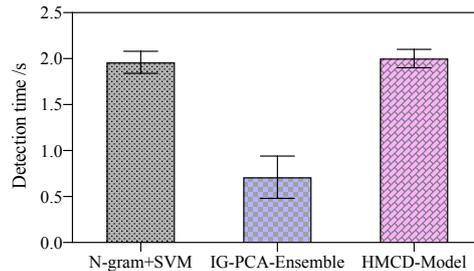


Figure 10: Comparison on detection time cost (per 10,000 samples).

## 7. Discussion

### 1) How to select the hyper-parameters?

We determine some of the best hyper-parameters (the number of packets in a flow, the byte size in a packet, *etc.*) of HMCT-Model through experiments. The magnitude of the

combination of hyper-parameters is very high and we cannot evaluate all solutions. As a rule of thumb, we use the method of controlling variables to cover the best solution as much as possible and to achieve a trade-off between detection performance and time cost. For instance, increasing the number of packets in a flow will improve the detection accuracy of the model, but the convergence time and detection speed will be slower. Therefore, we set two packets in a flow considering the trade-off between detection time and accuracy.

### **2) What are the limitations of HMCD-Model?**

HMCD-Model will cause misjudgment in two cases: First, some packets (such as packets in CIC-IDS-2017) have less packet content and insufficient features. Fewer features cannot give the model enough information. Second, the traffic generated by malware is indistinguishable from benign data. For instance, some malicious request packets have only two lines. However, benign packets in practice also have similar characteristics, which leads to inaccurate final judgment results of the model.

GAN can be used to generate HTTP-based malicious communication traffic, but it is generated based on existing frameworks. It still has certain limitations in many unknown attack scenes.

### **3) Why choose WGAN-GP to generate GAFs?**

Compared with vanilla GAN, WGAN-GP has stronger imitation ability and easier convergence, while maintaining the diversity of generated data [16]. HTTP-based application layer traffic belongs to text data, which is more difficult to generate than image data, and vanilla GAN is difficult to imitate. In this paper, when we try to generate GAFs with vanilla GAN, the model has a high probability of collapse and fails to converge. Therefore, WGAN-GP is relatively more suitable.

### **4) Can GAFs be used in other related methods to improve performance?**

In this paper, we utilize WGAN-GP to construct generated adversarial flows (GAFs) for data enhancement, to compensate for the insufficiency of real HTTP-based malicious communication traffic. The experiments in Section 6.3.2 show that GAFs can improve the generalization of HMCD-Model. We also try to use GAFs for other related methods. However, experimental results show that although GAFs can improve the performance of other methods, the improvement is not much, and sometimes even leads to performance degradation. There are two reasons here. First, other related methods have limitations in detecting HTTP-based malicious communication traffic and are difficult to comprehensively extract relevant features. Second, GAFs have a fixed composition, a flow consists of only one request and one response packet, and the represented feature space is limited.

Taking HMCT2020 (19-20) as an example, we use ep2 as the training set and test set, and increase GAFs when training N-gram+SVM and IG-PCA-Ensemble. The experimental results are shown in Tab 13. For N-gram+SVM, after adding GAFs in the training set, the precision rate is greatly improved (81.79% to 87.56%), but the recall rate is reduced (62.99% to 63.33%). For IG-PCA-Ensemble, after adding GAFs in the training set, the precision rate decreases (75.92% to 73.85%) and the recall rate increases (70.07% to 78.55%). In general, the improvement effect of GAFs on other methods is limited, and for HMCD-Model, the generalization can be improved comprehensively.

### **5) Can HMCD-Model be used for HTTPS encrypted traffic detection?**

Table 13: Generalization Comparison with GAFs in dataset HMCT-2020(19-20) (ep2 in Tab 7)

	HMCT-2020(19-20) (ep2)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
N-gram+SVM	81.79 <sup>+1.83</sup> <sub>-1.07</sub>	51.24 <sup>+1.31</sup> <sub>-1.24</sub>	62.99 <sup>+0.73</sup> <sub>-0.84</sub>	11.42 <sup>+1.08</sup> <sub>-1.34</sub>
N-gram+SVM (+GAFs)	87.56 <sup>+0.81</sup> <sub>-2.16</sub>	49.60 <sup>+2.09</sup> <sub>-1.82</sub>	63.33 <sup>+0.73</sup> <sub>-0.84</sub>	6.75 <sup>+1.04</sup> <sub>-0.98</sub>
IG-PCA-Ensemble	75.92 <sup>+6.05</sup> <sub>-6.21</sub>	70.07 <sup>+4.68</sup> <sub>-4.87</sub>	72.59 <sup>+0.60</sup> <sub>-0.45</sub>	23.02 <sup>+9.46</sup> <sub>-8.68</sub>
IG-PCA-Ensemble (+GAFs)	73.85 <sup>+5.45</sup> <sub>-6.52</sub>	78.55 <sup>+3.81</sup> <sub>-4.93</sub>	76.13 <sup>+0.74</sup> <sub>-0.58</sub>	24.96 <sup>+8.52</sup> <sub>-9.79</sub>
HMCD-Model	<b>97.94</b> <sup>+0.70</sup> <sub>-1.51</sub>	<b>99.39</b> <sup>+0.18</sup> <sub>-0.25</sub>	<b>98.66</b> <sup>+0.38</sup> <sub>-0.89</sub>	<b>2.10</b> <sup>+1.57</sup> <sub>-0.73</sub>

Currently, many studies have achieved good results in protocol classification and application identification [36, 11, 28] based on the strategy that directly analyzing encrypted traffic (such as HTTPS) without decryption. However, these methods are difficult to be effective for the refined detection of highly complex and covert attack traffic (such as HTTP-based malicious communication traffic). The current mainstream strategy is to deploy the detection system in the terminal instead of the middle. In this way, the natural decryption capability of the terminal can be used to convert HTTPS to HTTP. Many methods (Wang *et al.* [44], *etc.*) and systems (WAF [10], *etc.*) for malicious behavior detection are based on this strategy. We can decrypt HTTPS traffic into HTTP traffic and analyze the malicious behavior it contains. HMCD-Model based on the above strategy can be feasible as an end-to-end detection method. It can be applied to the IDSs of host/server.

## 8. Conclusion

For the detection of unknown HTTP-based malicious communication behavior, we propose HMCD-Model. We build a hybrid neural network to learn the hierarchical spatial-temporal features of traffic; use GAN to generate adversarial flows to make up for the lack of real traffic. In addition, we publish dataset HMCT-2020. Compared with the most representative methods at present, our model has obvious advantages in generalization, which can reach the  $F1 \approx 83.66\%$  in detection of real malware traffic. HMCD-Model can be applied to the detection of unknown HTTP-based malicious communication behavior and improve the capability of defenders against highly complex and covert network attack events.

In the future, we will further collect and expand datasets. In addition, we will fine-grain such malicious behavior scene to improve the ability of GAN-based malicious traffic generation.

## Acknowledgements

This work is supported by the National Key Research and Development Program of China (Grant No.2018YFB0804704), and the National Key Research and Development Program of China (Grant No.2019YFB1005201).

### A. Generalization Comparison with Classical Baseline Methods

We compared HMCD-Model with classical machine learning methods (Bayes, Decision Tree, SVM) and conventional neural networks (DNN, CNN, LSTM). We use data in HMCT-2020(18) as the training set and data in 82-Malware-Traffic as the test set. The experimental results are shown in Tab. 14. HMCD-Model is at least 5.53% better than other methods in terms of  $F1$ . Although Naive Bayes can achieve the  $R$  of 82.07%, the  $FPR$  is the highest, 28.03%. This is because Naive Bayes decreases the threshold for determining a sample as malicious, which means a lot of false positives. Similarly, the precision of SVM is  $P=99.87\%$  and the  $FPR$  is 0.08, but other metrics are very low. This is because SVM increases the threshold for determining a sample as malicious, which means a lot of false negatives. For instance, the  $F1$  of SVM is only 75.16%. Therefore, HMCD-Model has a better comprehensive detection performance than other classical methods.

Table 14: Comparison with classical baseline methods (ep4 in Tab 7)

	82-Malware-Traffic (ep4)			
	$P(\%)$	$R(\%)$	$F1(\%)$	$FPR(\%)$
Naive Bayes	74.55 <sup>+0.52</sup> <sub>-0.26</sub>	<b>82.07</b> <sup>+0.08</sup> <sub>-0.2</sub>	78.13 <sup>+0.32</sup> <sub>-0.12</sub>	28.03 <sup>+0.4</sup> <sub>-0.75</sub>
Decision Tree	99.7 <sup>+0.13</sup> <sub>-0.23</sub>	58.08 <sup>+1.77</sup> <sub>-4.06</sub>	73.37 <sup>+1.46</sup> <sub>-3.28</sub>	0.18 <sup>+0.14</sup> <sub>-0.08</sub>
SVM	<b>99.87</b> <sup>+0.08</sup> <sub>-0.08</sub>	60.26 <sup>+0.77</sup> <sub>-0.48</sub>	75.16 <sup>+0.57</sup> <sub>-0.35</sub>	<b>0.08</b> <sup>+0.05</sup> <sub>-0.05</sub>
DNN	99.83 <sup>+0.06</sup> <sub>-0.06</sub>	55.82 <sup>+0.59</sup> <sub>-0.5</sub>	71.6 <sup>+0.49</sup> <sub>-0.41</sub>	0.10 <sup>+0.03</sup> <sub>-0.04</sub>
CNN	94.95 <sup>+1.92</sup> <sub>-4.45</sub>	52.0 <sup>+2.37</sup> <sub>-2.16</sub>	67.15 <sup>+1.38</sup> <sub>-1.51</sub>	2.83 <sup>+2.87</sup> <sub>-1.14</sub>
LSTM	99.77 <sup>+0.11</sup> <sub>-0.06</sub>	55.44 <sup>+2.62</sup> <sub>-1.27</sub>	71.26 <sup>+2.13</sup> <sub>-1.03</sub>	0.13 <sup>+0.03</sup> <sub>-0.07</sub>
HMCD-Model	96.62 <sup>+2.26</sup> <sub>-4.27</sub>	73.77 <sup>+2.01</sup> <sub>-3.41</sub>	<b>83.66</b> <sup>+2.15</sup> <sub>-3.79</sub>	2.57 <sup>+3.26</sup> <sub>-1.71</sub>

## References

- [1] Aburomman, A.A., Reaz, M.B.I., 2017. A survey of intrusion detection systems based on ensemble and hybrid classifiers. *Computers & Security* 65, 135–152.
- [2] Altman, N.S., 1992. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician* 46, 175–185.
- [3] Cannady, J., 1998. Artificial neural networks for misuse detection, in: National information systems security conference, Baltimore. pp. 443–456.
- [4] Caviglione, L., Gaggero, M., Lalande, J.F., Mazurczyk, W., Urbański, M., 2015. Seeing the unseen: revealing mobile malware hidden communications via energy consumption and artificial intelligence. *IEEE Transactions on Information Forensics and Security* 11, 799–810.
- [5] Chandola, V., Banerjee, A., Kumar, V., 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 1–58.
- [6] Chen, E., Bates, T., 1996. Rfc1998: An application of the bgp community attribute in multi-home routing.
- [7] Cheng, A., 2019. Pac-gan: Packet generation of network traffic using generative adversarial networks, in: 2019 IEEE 10th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), IEEE. pp. 0728–0734.
- [8] Cheng, Q., Zhou, S., Shen, Y., Kong, D., Wu, C., 2021. Packet-level adversarial network traffic crafting using sequence generative adversarial networks. *arXiv preprint arXiv:2103.04794* .
- [9] Chowdhury, M.M.U., Hammond, F., Konowicz, G., Xin, C., Wu, H., Li, J., 2017. A few-shot deep learning approach for improved intrusion detection, in: 2017 IEEE 8th Annual Ubiquitous Computing, Electronics and Mobile Communication Conference (UEMCON), IEEE. pp. 456–462.
- [10] Clincy, V., Shahriar, H., 2018. Web application firewall: Network security models and configuration, in: 2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC), IEEE. pp. 835–836.
- [11] Di Martino, M., Quax, P., Lamotte, W., 2019. Realistically fingerprinting social media webpages in https traffic, in: Proceedings of the 14th International Conference on Availability, Reliability and Security, pp. 1–10.
- [12] Du, M., Li, F., Zheng, G., Srikumar, V., 2017. Deeplog: Anomaly detection and diagnosis from system logs through deep learning, in: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, pp. 1285–1298.
- [13] Du, Z., Ma, L., Li, H., Li, Q., Sun, G., Liu, Z., 2018. Network traffic anomaly detection based on wavelet analysis, in: 2018 IEEE 16th International Conference on Software Engineering Research, Management and Applications (SERA), IEEE. pp. 94–101.
- [14] Fukushima, K., Miyake, S., Ito, T., 1983. Neocognitron: A neural network model for a mechanism of visual pattern recognition. *IEEE transactions on systems, man, and cybernetics* , 826–834.
- [15] Ghafir, I., Svoboda, J., Prenosil, V., et al., 2015. A survey on botnet command and control traffic detection. *Int J Adv Comput Netw Secur* 5, 7580.
- [16] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C., 2017. Improved training of wasserstein gans, in: Advances in neural information processing systems, pp. 5767–5777.
- [17] Gupta, D., Khanna, A., SK, L., Shankar, K., Furtado, V., Rodrigues, J.J., 2019. Efficient artificial fish swarm based clustering approach on mobility aware energy-efficient for manet. *Transactions on Emerging Telecommunications Technologies* 30, e3524.
- [18] Hajisalem, V., Babaie, S., 2018. A hybrid intrusion detection system based on abc-afs algorithm for misuse and anomaly detection. *Computer Networks* 136, 37–50.
- [19] Hao, X., Jiang, Z., Xiao, Q., Wang, Q., Yao, Y., Liu, B., Liu, J., 2021. Producing more with less: A gan-based network attack detection approach for imbalanced data, in: 2021 IEEE 24th International Conference on Computer Supported Cooperative Work in Design (CSCWD), IEEE. pp. 384–390.
- [20] Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Computation* 9, 1735–1780.
- [21] Jan, S.T., Hao, Q., Hu, T., Pu, J., Oswal, S., Wang, G., Viswanath, B., 2020. Throwing darts in the

- dark? detecting bots with limited data using neural data augmentation, in: The 41st IEEE Symposium on Security and Privacy (IEEE SP).
- [22] Jose, S., Malathi, D., Reddy, B., Jayaseeli, D., 2018. A survey on anomaly based host intrusion detection system, in: Journal of Physics: Conference Series, IOP Publishing. p. 012049.
- [23] Karaboga, D., Basturk, B., 2008. On the performance of artificial bee colony (abc) algorithm. Applied soft computing 8, 687–697.
- [24] Kim, T., Kang, B., Rho, M., Sezer, S., Im, E.G., 2018. A multimodal deep learning method for android malware detection using various features. IEEE Transactions on Information Forensics and Security 14, 773–788.
- [25] Kwon, D., Kim, H., Kim, J., Suh, S.C., Kim, I., Kim, K.J., 2019. A survey of deep learning-based network anomaly detection. Cluster Computing , 1–13.
- [26] Li, J., Zhou, L., Li, H., Yan, L., Zhu, H., 2019. Dynamic traffic feature camouflaging via generative adversarial networks, in: 2019 IEEE Conference on Communications and Network Security (CNS), IEEE. pp. 268–276.
- [27] Lin, Z., Shi, Y., Xue, Z., 2018. Idsgan: Generative adversarial networks for attack generation against intrusion detection. arXiv preprint arXiv:1809.02077 .
- [28] Lotfollahi, M., Siavoshani, M.J., Zade, R.S.H., Saberian, M., 2020. Deep packet: A novel approach for encrypted traffic classification using deep learning. Soft Computing 24, 1999–2012.
- [29] Maki, M.C., Feller, W.J., 1989. Intrusion detection system. US Patent 4,879,544.
- [30] Moustafa, N., Slay, J., 2015. Unsw-nb15: a comprehensive data set for network intrusion detection systems (unsw-nb15 network data set), in: 2015 military communications and information systems conference (MilCIS), IEEE. pp. 1–6.
- [31] Ring, M., Schlör, D., Landes, D., Hotho, A., 2019. Flow-based network traffic generation using generative adversarial networks. Computers & Security 82, 156–172.
- [32] Ring, M., Wunderlich, S., Gründl, D., Landes, D., Hotho, A., 2017. Flow-based benchmark data sets for intrusion detection, in: Proceedings of the 16th European conference on cyber warfare and security, pp. 361–369.
- [33] Salo, F., Nassif, A.B., Essex, A., 2019. Dimensionality reduction with ig-pca and ensemble classifier for network intrusion detection. Computer Networks 148, 164–175.
- [34] Selvakumar, B., Muneeswaran, K., 2019. Firefly algorithm based feature selection for network intrusion detection. Computers & Security 81, 148–155.
- [35] Sharafaldin, I., Lashkari, A.H., Ghorbani, A.A., 2018. Toward generating a new intrusion detection dataset and intrusion traffic characterization., in: ICISSP, pp. 108–116.
- [36] Shen, M., Wei, M., Zhu, L., Wang, M., 2017. Classification of encrypted traffic with second-order markov chains and application attribute bigrams. IEEE Transactions on Information Forensics and Security 12, 1830–1843.
- [37] Shiravi, A., Shiravi, H., Tavallaee, M., Ghorbani, A.A., 2012. Toward developing a systematic approach to generate benchmark datasets for intrusion detection. computers & security 31, 357–374.
- [38] Shrestha, P.L., Hempel, M., Rezaei, F., Sharif, H., 2015. A support vector machine-based framework for detection of covert timing channels. IEEE Transactions on Dependable and Secure Computing 13, 274–283.
- [39] Song, J., Takakura, H., Okabe, Y., Eto, M., Inoue, D., Nakao, K., 2011. Statistical analysis of honeypot data and building of kyoto 2006+ dataset for nids evaluation, in: Proceedings of the first workshop on building analysis datasets and gathering experience returns for security, pp. 29–36.
- [40] Stolfo, S., et al., 1999. Kdd-99 dataset. Available on <http://www.kdd.ics.uci.edu/databases/kddcup99/kddcup99.html> kddcup99.html .
- [41] Sundermeyer, M., Schlüter, R., Ney, H., 2012. Lstm neural networks for language modeling, in: Thirteenth annual conference of the international speech communication association.
- [42] Suykens, J.A., Vandewalle, J., 1999. Least squares support vector machine classifiers. Neural processing letters 9, 293–300.
- [43] Tavallaee, M., Bagheri, E., Lu, W., Ghorbani, A.A., 2009. A detailed analysis of the kdd cup 99 data

- set, in: 2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications, IEEE. pp. 1–6.
- [44] Wang, S., Yan, Q., Chen, Z., Yang, B., Zhao, C., Conti, M., 2017. Detecting android malware leveraging text semantics of network flows. *IEEE Transactions on Information Forensics and Security* 13, 1096–1109.
- [45] Wang, W., Shang, Y., He, Y., Li, Y., Liu, J., 2020. Botmark: Automated botnet detection with hybrid analysis of flow-based and graph-based traffic behaviors. *Information Sciences* 511, 284–296.
- [46] White, B., 1963. *Principles of neurodynamics: Perceptrons and the theory of brain mechanisms*.
- [47] Yu, L., Zhang, W., Wang, J., Yu, Y., 2017. Seqgan: Sequence generative adversarial nets with policy gradient, in: *Proceedings of the AAAI conference on artificial intelligence*.
- [48] Zhou, Y., Cheng, G., Jiang, S., Dai, M., 2020. Building an efficient intrusion detection system based on feature selection and ensemble classifier. *Computer Networks* , 107247.
- [49] Zingo, P., Novocin, A., 2020. Can gan-generated network traffic be used to train traffic anomaly classifiers?, in: *2020 11th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*, IEEE. pp. 0540–0545.