



A multi-layer composite identification scheme of cryptographic algorithm based on hybrid random forest and logistic regression model

Ke Yuan^{1,2} · Yabing Huang^{1,3} · Zhanfei Du¹ · Jiabao Li^{4,5} · Chunfu Jia^{1,6}

Received: 11 December 2022 / Accepted: 3 August 2023 / Published online: 21 August 2023
© The Author(s) 2023

Abstract

Cryptographic technology can effectively defend against malicious attackers to attack sensitive and private information. The core of cryptographic technology is cryptographic algorithm, and the **cryptographic algorithm identification is the premise of in-depth analysis of cryptography**. In the cryptanalysis of unknown cryptographic algorithm, the primary task is to identify the cryptographic algorithm used in the encryption and then carry out targeted analysis. With the rapid growth of Internet data, the increasing complexity of communication environment, and the increasing number of cryptographic algorithms, the single-layer identification scheme of cryptographic algorithm faces great challenges in terms of **identification ability and stability**. To solve these problems, on the basis of existing identification schemes, this paper proposes a new cluster division scheme **CMSSBAM-cluster**, and then proposes a multi-layer composite identification scheme of cryptographic algorithm using a composite structure. The scheme adopts the method of cluster division and single division to identify various cryptographic algorithms. Based on the idea of ensemble, **the scheme uses the hybrid random forest and logistic regression (HRFLR) model** for training, and conducts research on a data set consisting of **1700 ciphertext files encrypted by 17 cryptographic algorithms**. In addition, two ensemble learning models, hybrid gradient boosting decision tree and logistic regression (HGBDTLR) model and hybrid k-neighbors and random forest (HKNNRF) model are used as controls to conduct controlled experiments in this paper. The experimental results show that multi-layer composite identification scheme of cryptographic algorithm based on HRFLR model has an accuracy rate close to **100% in the cluster division stage**, and the identification results are higher than those of the other two models in both the cluster division and single division stages. In the last layer of cluster division, the identification accuracy of ECB and CBC encryption modes in block cryptosystem is significantly higher than that of the other two classification models by 35.2% and 36.1%. In single division, the identification accuracy is higher than HGBDTLR with a maximum of 9.8%, and higher than HKNNRF with a maximum of 7.5%. At the same time, the scheme proposed in this paper has significantly improved the identification effect compared with the single division identification accuracy of 17 cryptosystem directly and the 17 classification accuracy of 5.9% compared with random classification, which indicates that multi-layer composite identification scheme of cryptographic algorithm based on HRFLR model has significant advantages in the accuracy of identifying multiple cryptographic algorithms.

Keywords Cryptanalysis · Cryptographic algorithm identification · Ensemble learning · Hybrid random forest and logistic regression · Cluster division identification · Single-layer identification

✉ Jiabao Li
li.jiabao@connect.um.edu.mo

¹ School of Computer and Information Engineering, Henan University, Kaifeng 475004, China

² Henan Province Engineering Research Center of Spatial Information Processing, Henan University, Kaifeng 475004, China

³ Henan Rural Credit Union, Zhengzhou 450018, China

⁴ Faculty of Science and Technology, University of Macao, Macao 519000, China

⁵ Henan Yinzu Security Technology, Zhengzhou 450003, China

⁶ College of Cybersecurity, Nankai University, Tianjin 300350, China

Introduction

At present, most cryptanalysis techniques are used to identify one or several specific cryptographic algorithms under the condition of known cryptography. However, in practical application scenarios, researchers usually cannot directly obtain the specific cryptographic system and cryptographic algorithm used in the generation of ciphertext, and it is difficult to predict and standardize the obtained ciphertext information within the scope of several fixed cryptographic systems. Therefore, it is a prerequisite for cryptanalysis to correctly identify the specific algorithm in the ciphertext cryptosystem. In addition, with the increasing complexity of network data and the increasing number of various cryptosystems and cryptographic algorithms in each category, how to design an identification scheme for a specific cryptographic algorithm in a multi-cryptosystem scenario has become an urgent problem to be solved, and continues to become a research hotspot.

The researchers' research on cryptosystem identification is divided into two stages, namely, the identification of classical cryptosystem and the identification of modern cryptosystem [1]. The two processing techniques used in classical ciphers are substitution and permutation. However, after ciphertext is encrypted by classical cipher, it often cannot achieve complete randomness, and will contain statistical laws inherent in plaintext or plaintext characters themselves. Compared with classical ciphers, modern ciphers have great advantages, and the ciphers often have high randomness. At present, the research methods of cryptographic algorithm identification scheme mainly include statistical method and machine learning method. With the continuous development of cryptography, the cryptographic identification schemes based on statistics can not meet researchers' needs by extracting various statistical indicators of ciphertext and comparing them.

Meanwhile, with the rapid development of artificial intelligence technology, it is increasingly being used to solve security issues. Currently, there is a close integration between artificial intelligence technology and security technology. It is not only reflected in popular fields such as smart cities [2], industrial Internet of Things [3], intelligent unmanned ground vehicles [4], and medical information systems [5, 6], but also plays an important role in the field of cryptography, especially in the direction of cryptographic algorithm identification.

Therefore, the method based on machine learning is introduced, and the classification method in machine learning is used as the main method of cryptosystem identification. The features and labels extracted from the ciphertext are put into the classifier, so as to train the classifier model, and then, the unknown ciphertext data are classified and identified on the trained classifier model. One of the main problems in

machine learning research is generalization, which refers to the ability to predict unknown data using learned models. In many cases, the theoretical and empirical performance of a single classifier is not as good as that of an ensemble model composed of several classifiers [7]. The essence of ensemble learning is to combine several weak classifiers to get a strong learner when dealing with learning tasks, and then use the strong learner to solve problems, so that the generalization of learning system can be effectively improved.

Most of the existing cryptographic algorithm identification schemes use single-layer identification, that is, to directly identify specific cryptographic algorithms. However, with the rapid growth of Internet data, the increasing complexity of communication environment, and the increasing number of cryptographic algorithms, the single-layer identification scheme of cryptographic algorithm faces great challenges in terms of identification ability and stability. Therefore, to solve the above problems, the cryptographic algorithm multi-layer identification scheme is proposed.

In the cryptographic algorithm multi-layer identification scheme, the first step is to identify the ciphertext by cluster division, and identify the cryptosystem category of the cryptographic algorithm used in ciphertext encryption. Then, under the specific cluster classification, the specific encryption algorithm used in ciphertext encryption is further identified. Although researchers have made great progress in designing multi-layer identification schemes for cryptographic algorithms, there are still many tasks and problems that need to be further studied. There are various classification methods for different cryptographic algorithms, and there is no unified evaluation standard. Therefore, the existing theoretical framework of the cryptographic algorithm multi-layer identification scheme also needs to be further improved.

According to previous research, we find that compared with the multi-classification identification of cryptographic algorithms, the binary identification of cryptographic algorithms tends to have higher accuracy and stronger stability. Therefore, this paper proposes a new clustering method CMSSBAM-cluster, which improves the theoretical system of multi-layer identification of cryptosystems, including a variety of specific encryption algorithms under classical cryptosystem, stream cryptosystem, block cryptosystem and asymmetric cryptosystem. Based on the idea of ensemble learning, combined with the single division identification scheme of cryptographic algorithm based on HRFLR model [8] proposed by Yuan et al. in 2023, a multi-layer composite identification scheme of cryptographic algorithm based on HRFLR is proposed to design a composite structure containing cryptosystem cluster division and specific cryptosystem single division. The experimental results show that the new proposed cluster division identification scheme has great advantages in accuracy and stability.

In “[Related work](#)”, a comprehensive review of related works is presented, summarizing the main findings, methods, and limitations of previous research. “[Overview of cryptographic algorithm identification](#)” provides an overview of cryptographic algorithm identification, focusing on methods for single-layer identification of cryptographic algorithms, multi-layer identification of cryptographic algorithms, and ciphertext feature extraction. “[Identification scheme](#)” presents the main idea and workflow of the CMSSBAM-cluster scheme and the multi-layer composite identification scheme of cryptographic algorithms based on the HRFLR model. In “[Experimental design and analysis](#)”, we describe the dataset, experimental results, and analysis. Finally, in “[Summary and prospect](#)”, we summarize the current work, highlight the advantages and improvements of the proposed approach compared to the existing methods, and outline future research directions.

Related work

In 2006, Dileep and Sekhar [9] proposed block cipher algorithm identification scheme based on support vector machine for AES, DES, 3DES, RC5, and Blowfish. In 2011, Manjula and Anitha [10] proposed a cryptographic algorithm identification scheme based on C4.5 Decision Tree (DT) for 11 cryptographic algorithms including Permutation, Substitution, DES, 3DES, AES, RC2, RC4, Blowfish, IDEA, RSA, and ECC. In 2012, Chou et al. [11] proposed to use Support Vector Machine (SVM) to identify cryptographic algorithms. By extracting 12 kinds of ciphertext features, the scheme conducts experiments on two-by-two classification and identification of block cipher algorithms in ECB and CBC modes. In 2013, DeSouza et al. [12] designed an identification scheme based on a neural network for cryptographic algorithms, such as MARS, RC6, and Twofish. In 2018, Mello and JAM [13] analyzed ciphertext files encrypted by seven encryption algorithms of ARC4, Blowfish, DES, Rijndael, RSA, Serpent, and Twofish in both ECB and CBC modes. The experimental results show that almost all cryptographic algorithms can be fully identified in ECB mode. In 2019, Zhao et al. [14] proposed a block cipher algorithm identification scheme based on randomness test, redesigned the ciphertext feature extraction method according to the National Institute of Standards and Technology (NIST) randomness test standard, and successfully completed the pairwise identification of six cryptographic algorithms. In the same year, Arvind and Ram [15] proposed a method using bit-plane image features and fuzzy decision criteria to realize the identification and isolation of images encrypted with the same key. In 2021, Ji et al. [16] proposed a SM4 block cryptosystem identification method based on randomness test. It is used to distinguish SM4 from other international standard block cipher algo-

rithms. The features of the ciphertext are extracted through the randomness test, and then, the extracted feature vectors are identified using machine learning algorithms. The results show that the accuracy rate is as high as 90%. In 2022, Grari et al. [17] proposed a cryptanalysis model based on deep neural network. The neural network takes plaintexts and their corresponding ciphertexts as inputs to predict the secret key of the cipher. The cryptanalysis problem is defined as a multi-label classification problem, and a multi-layer perceptron is used for prediction. The results show that it is better to treat the cryptanalysis problem as a multi-label classification problem. The above shows the related works of cryptographic algorithm single-layer identification scheme, as shown in Table 1.

Most of the existing cryptographic algorithm identification schemes use single-layer identification, that is, to directly identify specific cryptographic algorithms. However, with the continuous complexity of network data, it is difficult to directly restrict the obtained ciphertext data to several fixed cryptographic algorithms. Therefore, how to find a still applicable identification scheme in the presence of multiple cryptosystems is an important premise for ciphertext analysis, which makes how to design a multi-layer identification scheme for cryptographic algorithms a research hotspot. Nowadays, some researchers have designed and studied the multi-layer identification of cryptographic algorithms. In 2013, Mishra et al. [18] proposed a two-stage joint identification method for ciphertext data of different file sizes encrypted by AES, DES, and Blowfish. The scheme combined entropy feature, block length analysis, and dictionary analysis to identify them, which makes them achieve an average identification accuracy of 80%. In 2015, Wu Yang et al. [19] conducted a randomness measurement value distribution for the ciphertexts encrypted and generated by five cipher algorithms: AES, Camellia, DES, 3DES, and SMS4. And they proposed a multi-layer identification scheme; it used the K-means clustering algorithm for cluster division, and then classified and identified according to the cluster division results. The accuracy was close to 90%. In 2018, Huang et al. [20] proposed a cryptographic algorithms multi-layer identification scheme based on random forest, and initially gave a complete definition system for cryptosystem identification problems. The scheme introduced three cluster division methods: CM-cluster, CSN-cluster, and CSBP-cluster. The experimental results showed that compared with the single-layer identification scheme, the identification accuracy of the multi-layer identification scheme had better advantages. In 2020, Wang Xu et al. [21] used a staged identification scheme to identify multiple cryptosystems, and studied the influence of different ciphertext features on the effect of the identification scheme. Combined with the Relief system selection algorithm and heterogeneous ensemble learning, they further proposed a dynamic identification scheme that could adapt to a variety of cryptosystem identification scenarios,

Table 1 Related works of cryptographic algorithm single-layer identification scheme

Year	Authors	Method	Identification Target	Results
2006	Dileep and Sekhar	SVM	AES, DES, 3DES, RC5, blowfish	Successfully identified 5 cryptographic algorithms
2011	Manjula and Anitha	C4.5	Permutation, substitution, DES, 3DES, AES, RC2, RC4, blowfish, IDEA, RSA, ECC	Successfully identified 11 cryptographic algorithms
2012	Chou et al	SVM	MARS, RC6, Rijndael, Serpent, Twofish	Successfully classified and identified in ECB and CBC modes by extracting 12 ciphertext features
2013	DeSouza et al	NN	MARS, RC6, Twofish	Successfully identified cryptographic algorithms such as MARS, RC6 and Twofish
2018	Mello and JAM	C4.5, FT, PART, Complement Naive Bayes	ARC4, Blowfish, DES, Rijndael, RSA, Serpent, Twofish	Almost all cryptographic algorithms can be fully identified in ECB mode
2019	Zhao et al	Randomness test	AES, DES, 3DES, IDEA, Blowfish, Camellia	Successfully carried out the pairing identification of 6 cryptographic algorithms
2019	Arvind and Ram	Bit-plane image features, fuzzy decision	Images encrypted with the same key	Successfully achieved identification and isolation of images encrypted with the same key
2021	Ji et al	Randomness test, C4.5	SM4, AES, DES, 3DES, Blowfish, Cast	90% accuracy in distinguishing SM4 from other standard block cipher algorithms
2022	Grari et al	DNN	AES	Successfully formulated the cryptanalysis problem as a multi-label classification problem and used a multi-layer perceptron to make predictions

Table 2 Related works of cryptographic algorithm multi-layer identification scheme

Year	Authors	Method	Target algorithms	Results
2013	Mishra et al	Entropy, Block length analysis, dictionary analysis	AES, DES, blowfish	Average identification rate of 80%
2015	Wu Yang et al	Entropy, K-means clustering	AES, Camellia, DES, 3DES, SM4	Accuracy rate close to 90%
2018	Huang et al	Random forest, clustering methods	CM-cluster, CSN-cluster, CSBP-cluster	Significant improvement in identification accuracy compared to single-layer identification
2020	Wang Xu et al	Feature selection, heterogeneous ensemble learning	CM-cluster, CSN-cluster, CSBP-cluster	The dynamic identification scheme can adapt to various cryptosystem identification scenarios
2023	Zhao et al	Hamming weight distribution, XGB-LGBM ensemble learning	AES-128, AES-256, 3DES, DES, IDEA, blowfish, SM4, Cast, RC2, Camellia	Overall accuracy of 89.65%

and proved the feasibility of the scheme. In 2023, Zhao et al. [22] proposed a block cryptosystem identification scheme based on Hamming weight distribution. They designed a feature extraction method using Hamming weight distribution and an XGB-LGBM ensemble learning model with a multi-layer fusion structure. The above shows the related works of cryptographic algorithm multi-layer identification scheme, as shown in Table 2. The experiment carried out mixed identification of ten common block cipher algorithm, and the overall accuracy reached 89.65%.

Overview of cryptographic algorithm identification

Key notations

For presentation convenience, we list the key notations used in our work in Table 3.

The single-layer identification scheme of cryptographic algorithm

The core of the cryptographic algorithm identification scheme is to **find out the small differences in the spatial distribution of the ciphertext data to be identified**, and then distinguish the cryptographic algorithm to which the ciphertext data belong. In the identification work of cryptography algorithm, most identification schemes **adopt supervised learning mode**, that is, feature is regarded as a group of attributes, an identification task is regarded as a classification task, classifier model is trained on the training data set containing feature and algorithm labels, and then, the trained model is used to identify the testing set data. Although this scheme is straightforward and easy to implement, it puts the problem of cryptographic algorithm identification in the framework of general pattern identification, which cannot achieve the expected results in the face of the specificity of cryptographic algorithm identification and brings obstacles to technological innovation.

Reference [20] standardizes the identification of specific cryptographic algorithms under the cryptosystem, and puts the identification methods adopted in the scheme into the proposed theoretical framework, which promotes the depth and breadth of research. And it gives a preliminary definition system for the cryptosystem identification problem. This section integrates the basic elements of the cryptographic algorithm identification problem, and supplements and improves the relevant definitions.

Definition 1 (*Cryptographic algorithm identification*) Suppose there is a set of cryptographic algorithms $A = \{a_1, a_2, \dots, a_N\}$, where N represents the number of cryptographic algorithm. In the ciphertext-only scenario, for any

Table 3 Key notations used in this paper

Notation	Interpretation
HRFLR	Hybrid random forest and logistic regression
HGBDTLR	Hybrid gradient boosting decision tree and logistic regression
HKNNRF	Hybrid K-neighbors and random forest
ECB	Electronic codebook
CBC	Cipher block chaining
AES	Advanced encryption standard
DES	Data encryption standard
3DES	Triple data encryption algorithm
RC	Rivest cipher
DT	Decision tree
IDEA	International data encryption algorithm
RSA	Rivest–Shamir–Adleman
ECC	Ellipse curve cryptography
SVM	Support vector machine
NIST	National Institute of Standards and Technology
PC	Personal computer
MCICA	Multilayer composite identification of cryptographic algorithm
AI	Artificial intelligence
A	Cipher
F	Cipher text
J	Identification scheme
C	Cluster
fea	Ciphertext features
Lab	Label
CPO	Identification of execution encryption algorithm set
RA	Identification scheme
ξ	Random variable
U	Identification scheme
As	Identification scheme
\rightarrow	Mapping
Φ	Encryption algorithm set (quintuple)
Δ	A five tuple composed of elements
U	Identification scheme under single cryptographic algorithm clustering recognition
∂	Triplet

given cryptographic algorithm $a_i \in A$, $1 \leq i \leq N$, let F be the ciphertext file generated by encryption, and its ciphertext data is known. Suppose there is an identification scheme J , which can identify the cryptographic algorithm to which the encrypted ciphertext file F belongs when a_i is unknown, and the identification accuracy is h^J , then this process is called cryptographic algorithm identification, and it can be recorded as triple $\partial = (A, J, h^J)$.

Definition 2 (*The single-layer identification scheme of cryptographic algorithm*) In cryptographic algorithm identifica-

tion $\vartheta = (A, J, h^J)$, $\text{oper}_{\text{SLRP}}$ is the workflow in identifying a specific cryptographic algorithm, the ciphertext feature extracted from the ciphertext file F is denoted as fea , and RA is the identification algorithm adopted by the scheme. Then, the single-layer identification scheme of the cryptographic algorithm can be recorded as $J = (\text{oper}_{\text{SLRP}}, \text{fea}, \text{RA})$.

The specific $\text{oper}_{\text{SLRP}}$ consists of two stages, training and testing, and the specific process is as follows:

(1) The stage of training.

Step 1. Collect a set of n ciphertext files F_1, F_2, \dots, F_n , and the cryptographic algorithm used for each file is known;

Step 2. Extract ciphertext features from the ciphertext files and collect a set of ciphertext features $\text{FeaTr} = \{\text{feaTr}_i = \text{feaTr}_i^j \mid i = 1, 2, \dots, n, j = 1, 2, \dots, d\}$ wherein any ciphertext feature feaTr_i is a d -dimensional feature vector;

Step 3. Denote the cryptographic algorithm of the ciphertext files as an n -dimensional vector $\text{Lab} = (\text{lab}_1, \text{lab}_2, \dots, \text{lab}_n)$ as the label of the feature data, where n is the number of ciphertext files. The two-tuples $(\text{FeaTr}, \text{Lab})$ composed of the feature set FeaTr and the label set Lab is called the ciphertext feature set containing the cryptographic algorithm label;

Step 4. Use the two-tuples $(\text{FeaTr}, \text{Lab})$ as the input of the classifier RA to train the classification model.

(2) The stage of testing.

Step 1. Adopt the same feature extraction method to extract the feature of the file FT which is to be identified with unknown label of the cryptographic algorithm, and obtain the d -dimensional feature, denoted as FeaTe , $\text{FeaTe} = \{\text{feaTe}^j \mid j = 1, 2, \dots, d\}$.

Step 2. Input the feature data FeaTe into the trained classifier RA , and the classifier provides the cryptographic algorithm identification result of the ciphertext FT according to the feature data, that is, the cryptographic algorithm a_{FT} to which the ciphertext belongs. The above two stages of training and testing constitute a complete single-layer workflow of cryptographic algorithm. The single-layer identification of cryptographic algorithm used in this paper will be further elaborated in “A multi-layer composite identification scheme of cryptographic algorithm”. For example, in classical cryptosystem, the single classification model of specific cryptographic algorithms ξ_{HRFLR_C} .

The multi-layer identification scheme of cryptographic algorithm

The problem of the multi-layer identification scheme of cryptographic algorithm is generally divided into two parts. The

first is to identify the ciphertext cluster, that is, to identify the cryptosystem category to which the ciphertext encryption algorithm belongs. Then, the single division is performed in a specific “cluster” to identify a specific cryptographic algorithm. Based on previous research results, this section describes the cluster division, CMSSBAM-cluster, and the definition of single division-related cryptography algorithm under cluster division identification, so as to further improve the theoretical knowledge of layered identification of cryptography.

Definition of the cluster division

The section gives the general definition of cluster division, clarifies the meaning of cluster division identification of cryptosystem, and gives the definition of CMSSBAM-cluster.

Definition 3 (*The cluster division*) Suppose there is a set of cryptographic algorithms

$$A = \{a_1, a_2, \dots, a_N\}$$

And cluster (category) set

$$C = \{c_1, c_2, \dots, c_K\},$$

where N is the number of cryptographic algorithms, K is the number of cluster, and there is $K \leq N$. It is known that there is a surjective $f : A \rightarrow C$ from A to C . For any given cryptographic algorithm $a_i \in A$, it is assumed that the ciphertext data are known, and the ciphertext file F is generated by a_i encryption. The following definition of cluster division can be given: if a_i is unknown, there is a certain identification scheme O , and the mapping value $f(a_i)$ is identified with a certain identification accuracy rate q^0 . And record this process as a quintuple

$$\Theta = (A, C, f, O, q^0).$$

The mapping f is called cluster division mapping.

Definition 4 (*Cryptographic algorithm cluster division identification scheme*) In cryptographic algorithm cluster division $\Theta = (A, C, f, O, q^0)$, a cluster division identification scheme can be described by a triple, namely $O = (\text{CPO}, \text{fea}, \text{RA})$, where CPO is the workflow for performing cryptographic algorithm cluster identification, fea is the ciphertext features extracted from the ciphertext file, and RA is the identification algorithm adopted by the scheme.

Definition 5 (*CMSSBAM-Cluster*) There is a set of cryptographic algorithms $A = \{a_1, a_2, \dots, a_N\}$ and the set of clusters $c_{\text{CMSSBAM}} = \{c_C, c_M\}$, $c_M = \{c_A, c_{S_B}\}$, $c_{S_B} = \{c_{S_S}, c_{S_B}\}$, $c_{S_{B_M}} = \{c_{S_{B_{\text{ECB}}}}, c_{S_{B_{\text{CBC}}}}\}$ where $N \geq 4$. The cluster division mapping $f_{\text{CMSSBAM}} : A \rightarrow c_{\text{CMSSBAM}}$ satisfies

$$f_{\text{CMSSBA}}(A_i) = \begin{cases} c_C, & a_i \text{ is classical cryptosystem} \\ c_M, & a_i \text{ is modern cryptosystem} \end{cases} \begin{cases} c_A, & a_i \text{ is asymmetric cryptosystem} \\ c_{S_B}, & a_i \text{ is symmetric cryptosystem} \end{cases}$$

$$c_{S_B} \begin{cases} c_{S_S}, & a_i \text{ is stream cryptosystem} \\ c_{S_B}, & a_i \text{ is block cryptosystem} \end{cases} \begin{cases} c_{S_{B_{ECB}}}, & a_i \text{ is ECB mode} \\ c_{S_{B_{CBC}}}, & a_i \text{ is CBC mode.} \end{cases}$$

Let F be a ciphertext file generated by encryption of any given cryptographic algorithm a_i . If there exists a identification scheme O_{CMSSBAM} that can identify the mapping relation $f_{\text{CMSSBAM}}(a_i)$ with a identification accuracy of q_{CMSSBAM}^O when $a_i \in A$ is unknown, then the identification process is said to be CMSSBAM-cluster. Denote this process as

$$\Theta_{\text{CMSSBAM}} = (A, C_{\text{CMSSBAM}}, f_{\text{CMSSBAM}}, O_{\text{CMSSBAM}}, q_{\text{CMSSBAM}}^O).$$

Definition of the single division

The single division is the process of identifying a specific cryptographic algorithm within a cluster of a cryptosystem. Therefore, it is necessary to assume that there is a completed cluster division process before this in the single identification scenario.

Definition 6 (*The single division*) For a given cluster division process $\Theta = (A, C, f, O, q^O)$, there is a set of cryptographic algorithms $As = \{as_1, as_2, \dots, as_Z\} \subset A$, where z is the number of specific cryptographic algorithms, under the cluster division mapping f , if $f(as_i) = c_{As} \in C, \forall as_i \in As$ and $f^{-1}(c_{As}) = As$ are satisfied, then c_{As} is called the cluster to which the set of cryptographic algorithms As belongs. For any given cryptographic algorithm $as \in As$, let F be a ciphertext file generated by as encryption, if there exists an identification scheme U that identifies the ciphertext file F belonging to cryptographic algorithm as with accuracy q^U , then this process is called a single division of cryptographic algorithms under the cluster division identification, and we denote the quintuple consisting of these elements as $\Delta = (\Theta, As, c_{As}, U, p^U)$.

Definition 7 (*The single division identification scheme of cryptographic algorithm*) In the single division $\Delta = (\Theta, As, c_{As}, U, p^U)$ of the cryptographic algorithm, the triple $U = (\text{oper}, \text{fea}, \text{alg})$ can be used to represent the single division identification scheme of cryptographic algorithm under the cluster division identification. Among them, fea is the ciphertext feature, oper is the workflow of the single division identification scheme of cryptographic algorithm, alg is the identification algorithm used in the identification process.

It can be seen that when carrying out single division identification, it is necessary to clarify the specific clustering process Θ before that, the cryptographic algorithm subset As and the corresponding cluster c_{As} under a certain cryptosystem cluster category, and then concretize it.

Ciphertext feature extraction

In the cryptographic algorithm hierarchical identification problem, the extraction of ciphertext features is required in both the cluster division of the cryptographic and the single of the specific cryptographic algorithm, and the extracted ciphertext features serve as the input to the identification model and directly affect the identification results of the ciphertext. Therefore, the key to the identification task is whether the ciphertext features can be reasonably extracted, so that they can effectively portray its information characteristics.

The SP 800-22 standard [23–27] developed by the NIST has a comprehensive and targeted test method for all aspects of ciphertext sequences and sequence segmentation methods, which has a wide detection coverage and is favored by many researchers. In this section, NIST randomness detection [28, 29] theory is used as the theoretical basis for randomness testing with hypothesis testing, and in practical applications, the P-value method is commonly used to determine whether the source hypothesis is valid. And with reference to the existing ciphertext feature extraction methods [19, 30] based on randomness detection, 40 features are redesigned and collected to carry out a multi-layer composite identification task of cryptographic algorithms based on ensemble learning, in which the extraction of the cluster division features in CMSSBAM-cluster identification and the extraction of single division features under a certain cryptographic cluster use the same feature extraction method. The final identification of the cryptographic algorithm to which the ciphertext to be tested belongs is achieved through two stages of training and testing, respectively. Table 4 shows the 40 cluster (single) division features designed in this chapter based on the literature.

Table 4 List of 40 kinds of cluster (single) division features

Feature extraction method	Feature
The Runs Test	The_Runs_Test
The Longest Run Ones in A Block Test	The_longest_run_ones_in_a_block_test
The Binary Matrix Rank Test	The_binary_matrix_rank_test
The Non-Overlapping Template Matching Test	The_Non_Overlapping_Template_Matching_Test
The Maurers Universal Test	The_Maurers_Universal_Test
The Serial Test	The_Serial_Test_1
	The_Serial_Test_2
The Approximate Entropy Test	The_Approximate_Entropy_Test
The Cumulative Sums Test	The_Cumulative_Sums_Test_1
	The_Cumulative_Sums_Test_2
The Random Excursions Test	The_Random_Excursions_Test_1
	..
	The_Random_Excursions_Test_8
The Random Excursions Variant Test	The_Random_Excursions_Variant_Test_1
	..
	The_Random_Excursions_Variant_Test_18
The Overlapping Template Matching Test	The_Overlapping_Template_Matching_Test
The Linear Complexity Test	The_Linear_Complexity_Test
The Binary Matrix Rank Test	The_Binary_Matrix_Rank_Test
The Non-overlapping Template Matching Test	The_Non_Overlapping_Template_Matching_Test
The Maurers Universal Test	The_Maurers_Universal_Test

Identification scheme

The multi-layer identification scheme of cryptographic algorithms is composed of the upper layer cluster division identification scheme and the lower layer single identification scheme in a specific clustering of a certain cryptosystem. This section introduces the new proposed CMSSBAM-cryptosystem cluster division identification scheme, and proposes a multi-layer composite identification scheme based on CMSSBAM-cluster in the identification context including multiple cryptosystems.

CMSSBAM-cryptosystem cluster division identification scheme

This paper proposes a CMSSBAM-cryptosystem cluster division identification scheme based on Definitions 4 and def5, which can be written as $O = (CPO_{CMSSBAM}, \text{fea}, \text{HRFLR})$. As shown in Fig. 1, $CPO_{CMSSBAM}$ is the workflow of the scheme, including two stages of training and testing; fea is the feature extracted by the ciphertext feature extraction method based on NIST randomness detection; HRFLR is the hybrid random forest and logistic regression model identification algorithm. In the identification scheme, the cluster division basis is the cluster division label set $\text{Clust} = \{\text{clust}_i \mid i = 1, 2, \dots, 8\}$, which is classical cryptosystem, modern cryptosystem, asymmetric cryptosystem,

symmetric cryptosystem, stream cryptosystem, block cryptosystem, ECB mode, and CBC mode. $CPO_{CMSSBAM}$ has two stages of training and testing, as follows:

(1) The stage of training.

Step 1. Collect four ciphertext datasets containing labels. In the training stage, we need to train four cluster division models. When training the $x(x = 1, 2, 3, 4)$ cluster division model, we need to collect a set of ciphertext files $F(x)_1, F(x)_2, \dots, F(x)_n$ with cluster division labels $2x - 1$ and $2x$, where n is the number of files.

Step 2. Extract features from ciphertext files. A set of features $\text{FeaTr}(x) = \{\text{featr}(x)_i^j \mid i = 1, 2, \dots, n, j = 1, 2, \dots, d\}$ is obtained, where $\text{featr}(x)_i^j$ represents the j th feature of the i th ciphertext file.

Step 3. Each ciphertext file is used as a sample, and the cryptosystem cluster division labels of n samples are used as classification labels, denoted as $\text{CLab}(x) = \{\text{clab}(x)_i \mid i = 1, 2, \dots, n\}$. The two-tuples $(\text{FeaTr}(x), \text{CLab}(x))$ composed of $\text{FeaTr}(x)$ and $\text{CLab}(x)$ is recorded as the original data set $T(x)$, and each sample has d features.

Step 4. Submit the dataset $T(x)$ to the classification algorithm HRFLR. Then train the classification model

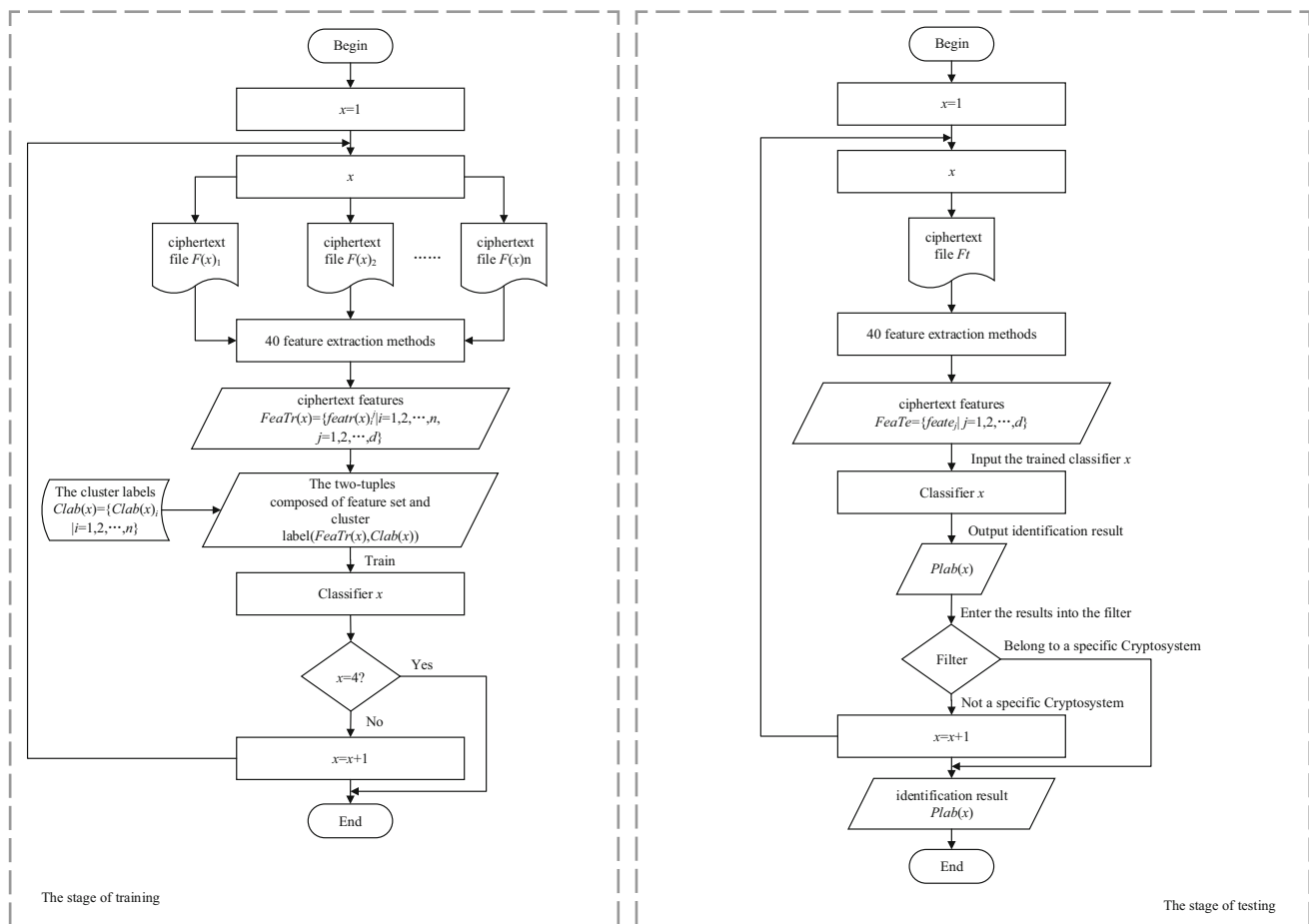


Fig. 1 The flow of the CMSSBAM-cryptosystem cluster division identification scheme

x , and finally get four cluster division classification models.

(2) The stage of testing.

Step 1. Collect and extract the ciphertext features of the ciphertext file Ft to be identified and denote them as $FeaTe = \{feate_j | j = 1, 2, \dots, d\}$.

Step 2. Input feature set $FeaTe$ to classifier 1, which gives the first cluster division identification result $PLab(1)$. Then enter $PLab(1)$ into the filter, and the filter is filtered according to $CLab(1)$.

Step 3. If $PLab(1)$ belongs to the classical cryptosystem, the classifier 1 directly outputs the cryptosystem cluster division label $PLab(1)$. Otherwise, record $PLab(1)$ and take the ciphertext feature $FeaTe$ as the input for the next step.

Step 4. Input feature set $FeaTe$ to classifier 2, which gives the second cluster division identification result $PLab(2)$. Then enter $PLab(2)$ into the filter, and the filter is filtered according to $CLab(2)$.

Step 5. If $PLab(2)$ belongs to asymmetric cryptosystem, the classifier 2 directly outputs the cryptosystem cluster division labels ($PLab(1)$, $PLab(2)$). Otherwise, record

$PLab(2)$ and take the ciphertext feature $FeaTe$ as the input for the next step.

Step 6. Input feature set $FeaTe$ to classifier 3, which gives the third cluster division identification result $PLab(3)$. Then enter $PLab(3)$ into the filter, and the filter is filtered according to $CLab(3)$.

Step 7. If $PLab(3)$ belongs to stream cryptosystem, the classifier 3 directly outputs the cryptosystem cluster division labels ($PLab(1)$, $PLab(2)$, $PLab(3)$). Otherwise, record $PLab(3)$ and take the ciphertext feature $FeaTe$ as the input for the next step.

Step 8. Input feature set $FeaTe$ to classifier 4, which gives the fourth cluster division identification result $PLab(4)$. Then input $PLab(4)$ into the filter, and the cryptosystem cluster division identification result of the ciphertext file Ft is ($PLab(1)$, $PLab(2)$, $PLab(3)$, $PLab(4)$).

The single division identification scheme

For the given single division identification setting $\Delta = (\Theta, As, c_{As}, U, p^U)$, the single division identification sche-

me adopted in this paper can be described as $U = (\text{SPO}, \text{fea}, \text{HRFLR})$ according to Definition 7. Among them, the SPO is a specific workflow, which is basically the same as the process described in $\text{oper}_{\text{SLRP}}$; fea is the feature of extracting ciphertext files; HRFLR is the hybrid random forest and logistic regression model identification algorithm.

A multi-layer composite identification scheme of cryptographic algorithm

This scheme is identified in the ciphertext-only scenario. Most of the existing cryptographic algorithm identification schemes are single-layer identification of cryptographic algorithm under a specific cryptographic system, but they are difficult in practical applications. Based on the previous

research results, this section further proposes a composite structure that includes a cryptosystem cluster division and a specific cryptosystem single division. In the cryptographic algorithm identification problem $\partial = (A, J, h^J)$, the multi-layer composite identification of cryptographic algorithm (MCICA) scheme based on CMSSBAM-cluster is described according to definition 2 as $J = (\text{oper}_{\text{MCICA}}, \text{fea}, \text{HRFLR})$. Among them, $\text{oper}_{\text{MCICA}}$ is the workflow of the multi-layer composite identification scheme of the cryptographic algorithm, fea represents the features collected by the scheme, and HRFLR represents the hybrid random forest and logistic regression algorithm adopted by identification.

The overall architecture of the scheme is shown in Fig. 2. The first cluster division is to distinguish between classical cryptosystem and modern cryptosystem, and then,

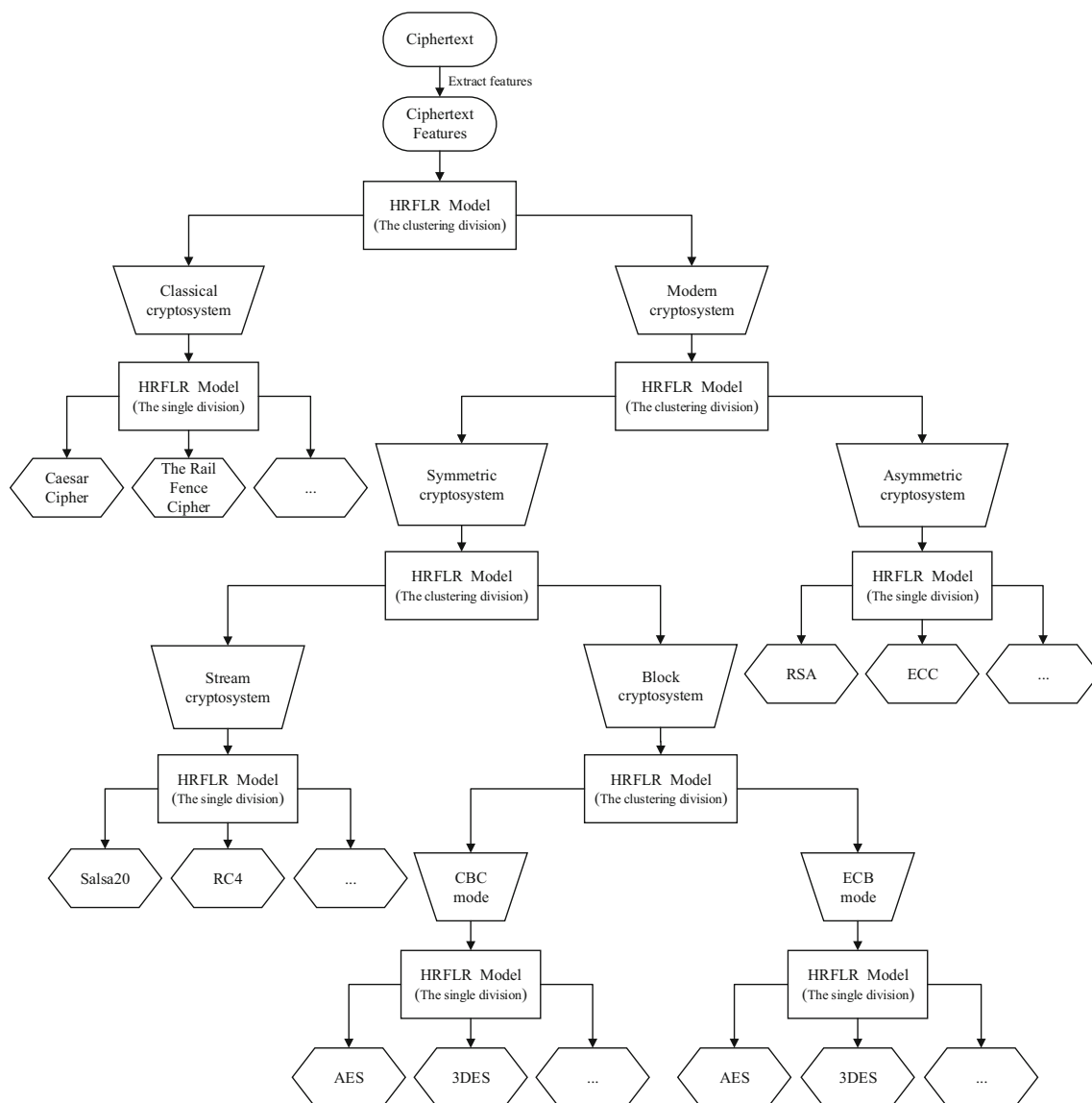


Fig. 2 Architecture of multi-layer composite identification scheme of cryptographic algorithm based on HRFLR model

single division is performed under the classical cryptosystem to complete specific classical cryptographic algorithm identification tasks. The second cluster division of modern cryptosystem is divided into symmetric cryptosystem and asymmetric cryptosystem, and then asymmetric cryptosystem is single-divided to complete the identification of specific asymmetric cryptographic algorithm. Symmetric cryptosystem is classified for the third cluster division. Symmetric cryptosystem is divided into stream cryptosystem and block cryptosystem. Then, stream cryptosystem is single-divided to complete the identification of specific stream cryptosystem. The block cryptosystem is classified into ECB encryption mode and CBC encryption mode for the fourth cluster division, and then, the cryptographic algorithm in each mode is identified separately, and finally, the specific cryptographic algorithm and encryption mode used for ciphertext encryption are identified.

The multi-layer composite identification scheme of cryptographic algorithm based on HRFLR model is composed of training stage and testing stage. The workflow of the scheme is mainly composed of the following steps:

Enter: The scheme requires the cryptosystem clustering category set $\text{clust}_1, \text{clust}_2, \dots, \text{clust}_8$, clustering category labels and cryptosystem labels are known for four groups of ciphertext file $\text{Ftr} = \{F(x)_j \mid x = 1, 2, 3, 4, j = 1, 2, \dots, n\}$, of which $F(x)_j$ is the j th ciphertext file involved in the x -time cluster division, and n is the total number of files. Finally, The scheme also needs to enter the testing ciphertext file Fte , whose cryptosystem is unknown.

Output: The classifier gives the identification result, that is, the cryptosystem a_{te} corresponding to the file to be tested.

(1) The stage of training.

Step 1: The scheme extracts the features of ciphertext files Ftr , and records the collected feature set as $\text{FeaTr}(x) = \{\text{featr}(x)_j^p \mid x = 1, 2, 3, 4, j = 1, 2, \dots, n, p = 1, 2, \dots, d\}$, where $\text{featr}(x)_j^p$ represents the p th feature of the j th ciphertext file in the x th clustering process.

Step 2. According to the cluster division identification scheme in “CMSSBAM-cryptosystem cluster division identification scheme”, four cluster division models are trained with ciphertext Ftr and feature set $\text{FeaTr}(x)$ as input.

Step 3. Select the ciphertext with the clustering labels of 1 in $F(1)_j$ to form the ciphertext set $F^c = \{F(1)_c \mid 1 \leq c \leq n\}$, and use the two-tuples composed of the ciphertext feature corresponding to F^c and the cryptographic algorithm labels $\text{Lab}^c = \{\text{lab}_1^c, \text{lab}_2^c, \dots, \text{lab}_n^c \mid 1 \leq c \leq n\}$ as input to train the single division classification model ξHRFLR_C of specific cryptographic algorithm in clas-

sical cryptosystem.

Step 4. Select the ciphertext with the clustering labels of 4 in $F(2)_j$ to form the ciphertext set $F^A = \{F(2)_a \mid 1 \leq a \leq n\}$, and use the two-tuples composed of the ciphertext feature corresponding to F^A and the cryptographic algorithm labels $\text{Lab}^A = \{\text{lab}_1^A, \text{lab}_2^A, \dots, \text{lab}_n^A \mid 1 \leq a \leq n\}$ as input to train the single division classification model ξHRFLR_A of specific cryptographic algorithm in asymmetric cryptosystem.

Step 5. Select the ciphertext with the clustering labels of 5 in $F(3)_j$ to form the ciphertext set $\text{FS} = \{F(3)_s \mid 1 \leq s \leq n\}$, and use the two-tuples composed of the ciphertext feature corresponding to FS and the cryptographic algorithm labels $\text{Lab}^{\text{FS}} = \{\text{lab}_1^{\text{FS}}, \text{lab}_2^{\text{FS}}, \dots, \text{lab}_n^{\text{FS}} \mid 1 \leq s \leq n\}$ as input to train the single division classification model ξHRFLR_S of specific cryptographic algorithm in stream cryptosystem.

Step 6. Select the ciphertext with the clustering labels of 7 in $F(4)_j$ to form the ciphertext set $F^{\text{ECB}} = \{F(4)_e \mid 1 \leq e \leq n\}$, and use the two-tuples composed of the ciphertext feature corresponding to F^{ECB} and the cryptographic algorithm labels $\text{Lab}^{\text{ECB}} = \{\text{lab}_1^{\text{ECB}}, \text{lab}_2^{\text{ECB}}, \dots, \text{lab}_n^{\text{ECB}} \mid 1 \leq e \leq n\}$ as input to train the single division classification model $\xi\text{HRFLR}_{\text{ECB}}$ of specific cryptographic algorithm under the ECB encryption mode of the block cryptosystem.

Step 7. Select the ciphertext with the clustering labels of 8 in $F(4)_j$ to form the ciphertext set $F^{\text{CBC}} = \{F(4)_b \mid 1 \leq b \leq n - e\}$, and use the two-tuples composed of the ciphertext feature corresponding to F^{CBC} and the cryptographic algorithm labels $\text{Lab}^{\text{CBC}} = \{\text{lab}_1^{\text{CBC}}, \text{lab}_2^{\text{CBC}}, \dots, \text{lab}_b^{\text{CBC}} \mid 1 \leq b \leq n - e\}$ as input to train the single division classification model $\xi\text{HRFLR}_{\text{CBC}}$ of specific cryptographic algorithm under the CBC encryption mode of the block cryptosystem.

(2) The stage of testing.

Step 1. Collect and extract the ciphertext feature of the Fte to be identified, and record it as $\text{FeaTe} = \{\text{feate}^j \mid j = 1, 2, \dots, d\}$. feate^j is the j th ciphertext feature of the file to be tested, and d is the feature dimension.

Step 2. Input the feature set FeaTe to the cluster division model ξHRFLR_1 , which gives the first cluster division identification result $\text{PLab}(1)$.

Step 3. Input $\text{PLab}(1)$ into the filter:

- (a) If $\text{PLab}(1)$ belongs to the classical cryptosystem, input FeaTe to ξHRFLR_C model to perform single identification of the cryptographic algorithm under the classical cryptosystem, and then give the identification result, that is, the cryptographic algorithm label a_{te} of the file to be tested;

- (b) If PLab(1) belongs to the modern cryptosystem, then input FeaTe for further cluster division.

Step 4. Input the feature set FeaTe to the cluster division model ξHRFLR_2 , which gives the second cluster division identification result PLab(2).

Step 5. Input PLab(2) into the filter:

- (a) If PLab(2) belongs to the symmetric cryptosystem, input FeaTe to ξHRFLR_A model to perform single identification of the cryptographic algorithm under the symmetric cryptosystem, and then give the identification result, that is, the cryptographic algorithm label a_{te} of the file to be tested;
- (b) If PLab(2) belongs to the symmetric cryptosystem, then input FeaTe for further cluster division.

Step 6. Input the feature set FeaTe to the cluster division model ξHRFLR_3 , which gives the third cluster division identification result PLab(3).

Step 7. Input PLab(3) into the filter:

- (a) If PLab(3) belongs to the stream cryptosystem, input FeaTe to ξHRFLR_S model to perform single identification of the cryptographic algorithm under the stream cryptosystem, and then give the identification result, that is, the cryptographic algorithm label a_{te} of the file to be tested;
- (b) If PLab(3) belongs to the block cryptosystem, then input FeaTe for further cluster division.

Step 8. Input the feature set FeaTe to the cluster division model ξHRFLR_4 , which gives the fourth cluster division identification result PLab(4).

Step 9. Input PLab(4) into the filter:

- (a) If PLab(4) belongs to the ECB encryption mode of block cryptosystem, input characteristic FeaTe to $\xi\text{HRFLR}_{\text{ECB}}$ model for single identification of the cryptographic algorithm in ECB mode, and then give the identification result, that is, the encryption algorithm label a_{te} of the file to be tested;
- (b) If PLab(4) belongs to the CBC encryption mode of block cryptosystem, input characteristic FeaTe to $\xi\text{HRFLR}_{\text{CBC}}$ model for single identification of the cryptographic algorithm in CBC mode, and then give the identification result, that is, the encryption algorithm label a_{te} of the file to be tested.

algorithm in multiple cryptosystems are collected. Specifically, this paper examines 17 cryptographic algorithms, one is AES, 3DES, CAST, Blowfish, and RC2 cryptographic algorithms in ECB and CBC modes of block cryptosystem; one is Caesar cipher and The rail fence cipher cryptographic algorithms under classical cryptosystem; one is Salsa20, RC4, and ChaCha20 cryptographic algorithms under the stream cryptosystem; the other is the RSA and ECC cryptographic algorithms under the asymmetric cryptosystem. The plaintexts used in the experiments in this paper are randomly generated by the uuid function of the python language. The experiments are conducted with 100 plaintext files of size 512 KB, and a total of 1700 ciphertext files are obtained by encrypting each of the above encryption algorithms. During the experiment, we set fixed encryption keys and initial variables for different block cipher algorithm. The encryption algorithm is implemented by the PyCryptodome package in the python Crypto library, and the HRFLR algorithm, HGB-DTLR algorithm, and HKNNRF algorithm are programmed by python language. The various cryptographic algorithms and parameter settings are shown in Table 5.

Evaluation criteria for classification results

In classification problems, accuracy, precision, and recall are commonly used as criteria to evaluate the effectiveness of scheme identification.

The correctness of the classification can be assessed by counting the number of correctly identified class examples (true positives), the number of correctly identified examples that do not belong to the class (true negatives), and examples that are incorrectly assigned to the class (false positives) or are not identified as class examples (false negatives) [31]. Let TP, TN, FP, and FN denote their corresponding sample numbers respectively, and the formulas (1)–(3) are used to calculate the accuracy rate, precision rate, and recall rate of the scheme respectively

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

Accuracy refers to the proportion of correct predictions made by the classification model, which is the ratio of true positives and true negatives to the total number of observations. Precision represents the proportion of true positives among all predicted positives by the model. Recall represents the proportion of true positives among all actual positives in the dataset.

Experimental design and analysis

Experimental subjects and datasets

The platform used for all experiments in this paper is a PC, and the ciphertext data generated by specific cryptographic

Table 5 Seventeen kinds of cryptographic algorithms and its parameter settings

Cryptographic algorithm	Key	Working mode	Way of implementation
AES	Fixed key	ECB	PyCryptodome
3DES	Fixed key	ECB	PyCryptodome
CAST	Fixed key	ECB	PyCryptodome
Blowfish	Fixed key	ECB	PyCryptodome
RC2	Fixed key	ECB	PyCryptodome
AES	Fixed key	CBC	PyCryptodome
3DES	Fixed key	CBC	PyCryptodome
CAST	Fixed key	CBC	PyCryptodome
Blowfish	Fixed key	CBC	PyCryptodome
RC2	Fixed key	CBC	PyCryptodome
Caesar cipher	Fixed key	–	PyCryptodome
The rail fence cipher	Fixed key	–	PyCryptodome
Salsa20	Fixed key	–	PyCryptodome
RC4	Fixed key	–	PyCryptodome
ChaCha20	Fixed key	–	PyCryptodome
RSA	Fixed key	–	PyCryptodome
ECC	Fixed key	–	PyCryptodome

Experimental results and analysis

Under the given cluster setting $\Theta_{\text{CMSSBAM}} = (A, C_{\text{CMSSBAM}}, f_{\text{CMSSBAM}}, O_{\text{CMSSBAM}}, q_{\text{CMSSBAM}}^O)$, we take into account the ciphertext data collected in the experimental data collection and processing stage, where A is a set composed of cryptographic algorithms under four cryptosystems, namely

$$A = \begin{cases} \text{CaesarCipher, Therailfencecipher} \\ \text{Salsa20, RC4, ChaCha20} \\ \text{AES, 3DES, CAST, Blowfish, RC2} \\ \text{RSA, ECC.} \end{cases}$$

Correspondingly, on the basis of the current setting A , we need to concretize the cluster map f_{CMSSBAM} , that is, let A_i represent the i th element of the set A , then

$$f_{\text{CMSSBA}}(A_i) = \begin{cases} c_C, 0 \leq i \leq 1 \\ c_M \begin{cases} c_{S_S}, 2 \leq i \leq 4 \\ c_{S_B}, 5 \leq i \leq 14 \\ c_A, 15 \leq i \leq 16. \end{cases} \end{cases}$$

For the cluster feature set $\text{FeaTr}(x) = \{\text{featr}(x)_j^p \mid x = 1, 2, 3, 4; j = 1, 2, \dots, n; p = 1, 2, \dots, d\}$, in the case of extracting the features of different clusters, we use the accuracy q_{CMSSBAM}^O as an indicator to confirm the feature data samples of the clusters, and perform repeated random sub-sampling verification. Eighty percent of the sample set is randomly selected for model training, and the remaining 20% of the samples are used as test data. In the same way, the same processing operation is performed for the single

division. In addition, when training the model, we divide the corresponding dataset according to the specific classification task that the current classifier is facing. The data set used by the classifier for training contains only the data with the label of the currently identified target.

Table 6 shows the best results of cluster division identification for three ensemble learning models based on HRFLR model, HGBDTLR model [32], and HKNNRF model [33] proposed by Yuan et al. for multi-layer composite identification scheme of cryptographic algorithm.

It can be seen from Table 6 that the first CM-cluster division is to identify classical and modern cryptosystem, and the identification accuracy of all three ensemble learning algorithms is 1.0, indicating that all test samples are correctly classified, and also indicating that classical and modern cryptosystem are very different and can be easily distinguished. The second SBA-cluster division is mainly designed to identify the symmetric cryptosystem and asymmetric cryptosystem in the modern cryptosystem, and the HRFLR identification accuracy is 0.992, which is 0.108 and 0.092 higher than the other two ensemble learning algorithms, respectively. The third SSB-cluster division is for the stream cryptosystem and the block cryptosystem in the symmetric cryptosystem, and the identification accuracy is 0.991, which is higher than the other two ensemble learning algorithms by 0.165 and 0.15, respectively. The fourth Mod-cluster division is to distinguish between ECB and CBC cryptographic patterns in the block cryptosystem, and the identification accuracy is 0.979, which is significantly higher than that of the other two classification models at 0.352 and 0.361. By comparing the precision and recall rates of

Table 6 Best identification results of multiple-cryptosystem cluster division

Classification model	Cluster method	Accuracy	Precision	Recall
HRFLR	CM-cluster	1.0	1.0	1.0
	SBA-cluster	0.992	0.992	0.992
	SSB-cluster	0.991	0.984	0.991
	Mod-cluster	0.979	0.961	0.979
HGBDTLR	CM-cluster	1.0	1.0	1.0
	SBA-cluster	0.884	0.804	0.884
	SSB-cluster	0.826	0.819	0.826
	Mod-cluster	0.627	0.627	0.627
HKNNRF	CM-cluster	1.0	1.0	1.0
	SBA-cluster	0.9	0.91	0.9
	SSB-cluster	0.841	0.868	0.841
	Mod-cluster	0.618	0.616	0.618

the three ensemble learning models in the four-layer cluster division separately, it can be found that the HRFLR model has the best identification results in all the same layer of cluster division methods, so the proposed multi-layer composite identification scheme of cryptographic algorithm in this paper is completed based on the HRFLR model.

From the cluster scheme and the accuracy of the cluster identification, the best results in Table 6 show that after the first cluster division, i.e., CM-cluster division, the two classical cipher algorithms, Caesar cipher and The rail fence cipher, are identified in a binary classification, and the remaining several modern cipher algorithms are further clustered. The second cluster division, namely SBA-cluster division, divides all samples into symmetric cryptosystem and asymmetric cryptosystem. Then, the asymmetric cryptosystem including RSA and ECC is identified by binary classification, and the samples belonging to the symmetric cryptosystem cluster are further clustered. The third cluster division, namely SSB-cluster division, divides the remaining samples into two cryptosystems, block cryptosystem and stream cryptosystem, and performs multi-classification and identification of three stream ciphers including Salsa20, RC4, and ChaCha20. Then, the samples belonging to the block cryptosystem are clustered in the next step. The fourth cluster division, Mod-cluster division, aims to distinguish the ECB and CBC encryption modes in the block cryptosystem. Then, the samples belonging to these two encryption modes are, respectively, identified by multiple classifications of cryptographic algorithms. In this way, the multi-layer composite identification of cryptographic algorithms based on HRFLR model is realized.

Table 7 shows the best identification results of multi-layer composite identification scheme of cryptographic algorithm based on the three ensemble learning models of HRFLR, HGBDTLR, and HKNNRF. It can be seen from Table 7 that the identification accuracy of HRFLR based on the first clus-

ter division, i.e., the Caesar cipher and the rail fence cipher under the classical cryptosystem and the modern cryptosystem, is 0.875, which is an improvement of 0.075 compared with the other two ensemble learning models. The HRFLR single division identification accuracy of the two asymmetric cipher algorithms, RSA and ECC, based on the second cluster division, i.e., symmetric cryptosystem and asymmetric cryptosystem, is 0.775, which is higher than that of the other two models by 0.098 and 0.033, respectively. Based on the third cluster division, i.e., cluster division of block ciphers and stream ciphers, the HRFLR single division identification accuracy of Salsa20, RC4, and ChaCha20 is 0.817, which is higher than that of the other two models at 0.045. Based on the fourth cluster division, i.e., the ECB and CBC modes of the block cryptosystem, the HRFLR single division accuracy of the five cryptographic algorithms belonging to the ECB mode is 0.242, and the HRFLR single division accuracy of the cryptographic algorithms belonging to the CBC mode is 0.239, which is higher than the identification accuracy of the other two models of 0.003, 0.01 and 0.017 and 0.011, respectively. By comparing the precision and recall of the three ensemble learning models in the five single division separately, it can be found that these two metrics of HRFLR are still higher than the other two models, further indicating that HRFLR is more applicable to the scheme proposed in this paper.

Figure 3 shows the scatter plot of the distribution of the ciphertext samples generated by the encryption of 17 cryptographic algorithms among the 40 ciphertext features extracted by randomness test. As shown in the figure, except that the two encryption algorithms belonging to the classical cryptosystem are greatly affected by the characteristics of The_Random_Excursions_Va-riant_Test series, various samples are mixed in the central area.

This shows that the two encryption algorithms belonging to the classical cryptosystem in the experiment are easier to

Table 7 Best identification results of the multi-layer composite identification of cryptographic algorithms

Method of use	Solution architecture	Working mode	Cryptographic algorithm	Accuracy	Precision	Recall
HRFLR	Layered	–	Caesar - The rail fence cipher	0.875	0.875	0.875
	Layered	–	RSA-ECC	0.775	0.777	0.775
	Layered	ECB	Five block cipher algorithms	0.242	0.234	0.242
	Layered	CBC	Five block cipher algorithms	0.239	0.24	0.239
	Layered	–	Salsa20-RC4-ChaCha20	0.817	0.821	0.817
	Single layer	–	All cryptographic algorithms	0.226	0.225	0.226
HGBDTLR	Layered	–	Caesar - The rail fence cipher	0.8	0.802	0.8
	Layered	–	RSA-ECC	0.677	0.677	0.677
	Layered	ECB	Five block cipher algorithms	0.239	0.238	0.239
	Layered	CBC	Five block cipher algorithms	0.222	0.221	0.222
	Layered	–	Salsa20-RC4-ChaCha20	0.772	0.791	0.772
	Single layer	–	All cryptographic algorithms	0.205	0.214	0.205
HKNNRF	Layered	–	Caesar - The rail fence cipher	0.8	0.812	0.8
	Layered	–	RSA-ECC	0.742	0.753	0.742
	Layered	ECB	Five block cipher algorithms	0.232	0.234	0.232
	Layered	CBC	Five block cipher algorithms	0.228	0.23	0.228
	Layered	–	Salsa20-RC4-ChaCha20	0.772	0.791	0.772
	Single layer	–	All cryptographic algorithms	0.208	0.22	0.208
Random classification	Single layer	–	All cryptographic algorithms	0.059	–	–

Fig. 3 Scatter plot of 17 cryptographic algorithms and feature distribution

identify, while the classification results of the modern cryptosystem have little correlation with the original dimension of the sample.

In this section, two ensemble learning models, HGBDTLR and HKNNRF, and three classification models based on 17 classification single division and random division identification are used as controls to conduct controlled experiments and compare the best experimental results. By comparing the accuracy, precision, and recall of the three ensemble learning models under four levels of cluster division and five single division cases, respectively, it is found that the HRFLR model used in the scheme of this paper is higher than the other

two in these three indexes, indicating that the effect of using the HRFLR model for multi-layer composite identification of cryptographic algorithm is better, and therefore, the proposed scheme in this paper is completed based on HRFLR. By comparing the single division identification accuracy of 17 cryptographic systems directly under the same ensemble learning model and the 17 division identification accuracy of 0.059 compared with random classification, all three ensemble learning models have a significant improvement in the identification effect of the multi-layer composite identification scheme of cryptographic algorithm, indicating that

the multi-layer composite identification scheme of cryptographic algorithm is feasible and effective.

Summary and prospect

In this paper, the problem of multi-layer identification of cryptographic algorithm under multi-cryptosystem is studied. Based on the existing theoretical framework of cryptographic algorithm identification, this paper improves the workflow of cryptographic algorithm identification, further proposes a new cluster identification scheme CMSSBAM-cluster identification scheme, and introduces the main idea and process of the scheme. On these basis, this paper adopts a composite structure and further proposes a multi-layer composite identification scheme, that is, a multi-layer composite identification of cryptographic algorithm scheme based on ensemble learning HRFLR model. In the case of multiple cryptosystems, this paper adopts the method of clustering and single-partitioning to identify 17 cryptographic algorithms, which cover classical cryptosystem, asymmetric cryptosystem, stream cryptosystem, block cryptosystem, and the ECB mode and CBC mode under the block cryptosystem. The experimental results show that multi-layer composite identification scheme of cryptographic algorithm based on HRFLR model has an accuracy rate close to 100% in the cluster division stage, and the identification results are higher than those of the other two models in both the cluster division and single division stages. In the last layer of cluster division, the identification accuracy of ECB and CBC encryption modes in block cryptosystem is significantly higher than that of the other two classification models by 35.2% and 36.1%. In single division, the identification accuracy is higher than HGBDTLR with a maximum of 9.8%, and higher than HKNNRF with a maximum of 7.5%. At the same time, the scheme proposed in this paper has significantly improved the identification effect compared with the single division identification accuracy of 17 cryptosystem directly and the 17 classification accuracy of 5.9% compared with random classification. It indicates that the multi-layer composite identification scheme of cryptographic algorithms proposed in this paper has obvious improvement significance in terms of theoretical support, scheme design, and identification accuracy. Therefore, in the future research on the identification of cryptographic algorithms, we can focus on the problem of multi-layer identification of cryptographic algorithm.

Although the proposed multi-layer composite identification scheme based on HRFLR model in this paper has achieved good experimental results, there is still room for improvement. With the explosion of deep learning, the latest deep learning algorithms have far exceeded the traditional machine learning algorithms for data prediction and clas-

sification accuracy. In addition, it is difficult to identify ciphertext files encrypted by ECB and CBC with the same cipher algorithm belonging to the block cryptosystem cluster. It is also difficult to identify different cryptographic algorithms using the same encryption mode in block cryptosystem cluster. In view of the above problems, the future research direction can focus on the use of deep learning models for multi-layer composite identification of cryptographic algorithms, and at the same time, for the encryption mechanism and structure of different cryptographic algorithms, the identification task can be carried out from the characteristics of the cryptographic algorithms themselves. In addition, the validity of cryptographic algorithm identification based on artificial intelligence (AI) technology also faces threats, such as spoofing and adversarial attacks, as well as interpretability, which are also issues to be further discussed in our future work. As a new idea, the cryptographic algorithm multi-layer identification based on AI technology scheme is worth exploring further, which has certain positive significance for future research on cryptographic algorithm identification.

Acknowledgements This work was supported by the National Key Research and Development Program (2018YFA0704703), the Fundamental Research Funds for the Central Universities of China, the National Natural Science Foundation of China (61972215, 61972073, 62172238), the Natural Science Foundation of Tianjin (20JCZDJC00640), the Key Specialized Research and Development Program of Henan Province (222102210062), the Basic Higher Educational Key Scientific Research Program of Henan Province (22A413004), and the National Innovation Training Program of University Student (202110475072).

Data availability The data can be made available on request from the corresponding author.

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Zhicheng Z (2018) The research of cryptosystem recognition scheme based on machine learning. Master's thesis,

- Support Forces Information Engineering University. CNKI:CDMD:2.1018.841883
2. Javed AR, Shahzad F, ur Rehman S, Zikria YB, Razzak I, Jalil Z, Xu G (2022) Future smart cities: requirements, emerging technologies, applications, challenges, and future aspects. *Cities* 129:103794. <https://doi.org/10.1016/j.cities.2022.103794>
 3. Alazab M, Gadekallu TR, Su C (2022) Guest editorial: security and privacy issues in industry 4.0 applications. *IEEE Trans Ind Inform* 18(9):6326–6329. <https://doi.org/10.1109/TII.2022.3164741>
 4. Han Z, Yang Y, Wang W, Zhou L, Gadekallu TR, Alazab M, Gope P, Su C (2022) Rssi map-based trajectory design for ugv against malicious radio source: a reinforcement learning approach. *IEEE Trans Intell Transp Syst* 24(4):4641–4650. <https://doi.org/10.1109/TITS.2022.3208245>
 5. El Zarif O, Haraty RA (2020) Toward information preservation in healthcare systems. In: Miltiadis D, Lytras MD, Sarirete A (eds) *Innovation in health informatics*. Academic Press, Next Gen Tech Driven Personalized Med & Smart Healthcare, pp 163–185. <https://doi.org/10.1016/B978-0-12-819043-2.00007-1>, <https://www.sciencedirect.com/science/article/pii/B9780128190432000071>. (ISBN:978-0-12-819043-2)
 6. Dhasarathan C, Hasan MK, Islam S, Abdullah S, Mokhtar UA, Javed AR, Goundar S (2023) COVID-19 health data analysis and personal data preserving: a homomorphic privacy enforcement approach. *Comput Commun* 199:87–97. <https://doi.org/10.1016/j.comcom.2022.12.004>
 7. Brown G, Wyatt J, Harris R, Yao X (2005) Diversity creation methods: a survey and categorisation. *Inf Fusion* 6(1):5–20. [https://doi.org/10.1016/S1566-2535\(04\)00037-5](https://doi.org/10.1016/S1566-2535(04)00037-5)
 8. Yuan K, Huang Y, Li J, Jia C, Yu D (2023) A block cipher algorithm identification scheme based on hybrid random forest and logistic regression model. *Neural Process Lett* 55:3185–3203. <https://doi.org/10.1007/s11063-022-11005-2>
 9. Dileep AD, Sekhar CC (2006) Identification of block ciphers using support vector machines. In: *The 2006 IEEE international joint conference on neural network proceedings*. IEEE, pp 2696–2701. <https://doi.org/10.1109/IJCNN.2006.247172>
 10. Manjula R, Anitha R (2011) Identification of encryption algorithm using decision tree. In: *International conference on computer science and information technology*. Springer, pp 237–246. https://doi.org/10.1007/978-3-642-17881-8_23
 11. Chou JW, Lin SD, Cheng CM (2012) On the effectiveness of using state-of-the-art machine learning techniques to launch cryptographic distinguishing attacks. In: *Proceedings of the 5th ACM Workshop on Security and Artificial Intelligence*, Raleigh, North Carolina, USA 6:105–110. <https://doi.org/10.1145/2381896.2381912>
 12. De Souza WA, Tomlinson A (2013) A distinguishing attack with a neural network. In: *2013 IEEE 13th international conference on data mining workshops*. IEEE, pp 154–161. <https://doi.org/10.1109/ICDMW.2013.116>
 13. de Mello FL, Xexeo JA (2018) Identifying encryption algorithms in ecb and cbc modes using computational intelligence. *J Univ Comput Sci* 24(1):25–42
 14. Zhao Z, Zhao Y, Liu F (2018) Recognition scheme of block cipher system based on randomness test. *J Cryptogr* 6(2):177–190. <https://doi.org/10.13868/j.cnki.jcr.000293>
 15. Ratan R et al (2020) Identifying traffic of same keys in cryptographic communications using fuzzy decision criteria and bit-plane measures. *Int J Syst Assur Eng Manag* 11(2):466–480. <https://doi.org/10.1007/s13198019-00878-7>
 16. Ji W, Li Y, Qin B (2021) Identification of sm4 block cipher system based on randomness characteristics. *Appl Res Comput*. <https://doi.org/10.19734/j.issn.1001-3695.2021.01.0019>
 17. Grari H, Zine-Dine K, Azouaoui A, Lamzabi S (2022) Deep learning-based cryptanalysis of a simplified aes cipher. *Int J Inf Secur Priv (IJISP)* 16(1):1–16. <https://doi.org/10.4018/IJISP.300325>
 18. Mishra S, Bhattacharjya A (2013) Pattern analysis of cipher text: a combined approach. In: *2013 International conference on recent trends in information technology (ICRTIT)*. IEEE, pp 393–398. <https://doi.org/10.1109/ICRTIT.2013.6844236>
 19. Wu Y, Wang T, Xing M, Li J (2015) Recognition scheme of block cipher algorithm based on distribution characteristics of random metric of ciphertext. *J Commun* 4:147–155. <https://doi.org/10.11959/j.issn.1000-436x.2015107>
 20. Huang L, Zhao Z, Zhao Y (2018) A two-stage cryptosystem recognition scheme based on random forest. *Chin J Comput* 41(2):382–399. <https://doi.org/10.11897/SPJ.1016.2018.00382>
 21. Wang X, Chen Y, Wang Q, Chen J (2021) Cryptosystem identification scheme combining feature selection and ensemble learning. *Comput Eng* 47(1):139–145. <https://doi.org/10.19678/j.issn.1000-3428.0056918>
 22. Zhao L, Chi Y, Xu Z, Yue Z (2023) Block cipher identification scheme based on hamming weight distribution. *IEEE Access* 11:21364–21373. <https://doi.org/10.1109/ACCESS.2023.3249753>
 23. Rukhin A, Soto J, Nechvatal J, Smid M, Barker E, Leigh S, Levenson M, Vangel M, Banks D, Heckert N, Dray J, Vo S (2001) *A statistical test suite for random and pseudorandom number generators for cryptographic applications*. Special Publication (NIST SP), National Institute of Standards and Technology, Gaithersburg, MD, Washington DC
 24. Minyong Q, Jinxin D, Quanke P (2008) Research and design of sequence randomness test system in information security. *Comput Eng Des* 29(6):1453–1455 (CNKI:SUN:SJSJ.0.2008-06-042)
 25. Liu J, Qu Q (2011) Randomness tests of several chaotic sequences. *Jisuanji Gongcheng yu Yingyong (Comput Eng Appl)* 47(5):46–49. <https://doi.org/10.3778/j.issn.1002-8331.2011.05.016>
 26. Shi G, Kang F, Gu H (2009) Research and implementation of randomness tests. *Comput Eng* 35(20):145–147. <https://doi.org/10.3969/j.issn.1000-3428.2009.20.051>
 27. Liu Z (2011) *The study of statistical tests of cryptographic algorithm*. Master's thesis, Xidian University. 10.7666/d.y1866917
 28. Mao S, Wang J, XL P (2006) *Advanced mathematical statistics*. Higher Education Press, Beijing
 29. Sheng Z (2001) *Probability theory and mathematical statistics*. Higher Education Press, Beijing
 30. Hongchao L (2018) *Research on cryptographic algorithm recognition based on ciphertext feature*. Master's thesis, Xidian University. CNKI:CDMD:2.1019.011746
 31. Sokolova M, Lapalme G (2009) A systematic analysis of performance measures for classification tasks. *Inf Process Manag* 45(4):427–437. <https://doi.org/10.1016/j.ipm.2009.03.002>
 32. Yuan K, Huang Y, Du Z, Li J, Jia C (2022) Block cipher algorithm identification scheme based on hybrid gradient boosting decision tree and logistic regression model. *Adv Eng Sci* 54(4):218–227. <https://doi.org/10.15961/j.jsuese.202100341>
 33. Yuan K, Yu D, Feng J, Yang L, Jia C, Huang Y (2022) A block cipher algorithm identification scheme based on hybrid k-nearest neighbor and random forest algorithm. *PeerJ Comput Sci* 8:1110. <https://doi.org/10.7717/peerj-cs.1110>