

Dear authors and members of the Computer Science Journal of Moldova team,

in the first issue of our journal (Volume 1, Number 1(1), 1993), I affirmed and expressed the hope:



*“Dear colleagues,
You have opened the first number of our journal. For us, it is a real symbol of our independence and we want it to be useful and interesting for you.*

...

We hope to accumulate here the main results on Computer Science obtained in Moldova and abroad and invite you to cooperate with us.”

Now we are at the last issue of Volume 30 of the journal (Volume 30, Number 3(90), 2022). We can note that along the way the journal had fruitful collaborations with colleagues from Europe, America, Asia, and Africa who published over 340 articles with us. I thank the authors, who considered our journal suitable to promote their results. And I cannot fail to mention the special contribution of the reviewers, whom I thank from the bottom of my heart.

The “Computer Science Journal of Moldova” is an Open Access journal, indexed in SCOPUS, Web of Science, EBSCO, Mathematical Reviews, Zentralblatt MATH, DOAJ, as well as in other Databases and journal indexing and abstracting services. We have an Editorial Board made up of outstanding personalities in Computer Science thanks to whom the goals announced in the journal’s first issue were achieved.

We hope that in the coming years we will be able to reach new achievements, make possible further progress, and have productive authors with excellent results.

I wish success to all the authors, reviewers, and the team of the Computer Science Journal of Moldova,

Constantin Gaidric, Editor-in-Chief

An Intelligent Detection of Malicious Intrusions in IoT Based on Machine Learning and Deep Learning Techniques

Saman Iftikhar, Danish Khan, Daniah Al-Madani,
Khattab M. Ali Alheeti, Kiran Fatima

Abstract

The devices of the Internet of Things (IoT) are facing various types of attacks, and IoT applications present unique and new protection challenges. These security challenges in IoT must be addressed to avoid any potential attacks. Malicious intrusions in IoT devices are considered one of the most aspects required for IoT users in modern applications. Machine learning techniques are widely used for intelligent detection of malicious intrusions in IoT. This paper proposes an intelligent detection method of malicious intrusions in IoT systems that leverages effective classification of benign and malicious attacks. An ensemble approach combined with various machine learning algorithms and a deep learning technique, is used to detect anomalies and other malicious activities in IoT. For the consideration of the detection of malicious intrusions and anomalies in IoT devices, UNSW-NB15 dataset is used as one of the latest IoT datasets. In this research, malicious and normal intrusions in IoT devices are classified with the use of various models.

Keywords: Malicious Intrusions, Anomaly detection, Machine Learning, Deep Learning, Classification, IoT dataset.

1 Introduction

The practice of integrating technologies such as sensors and software in everyday objects is on the rise. The idea is to enhance automation and

enable the transfer of data without the need for computer-human interaction. Devices such as coffee machines, stoves, dryers, washers, and refrigerators contain smart capabilities, which stretch their functionality and user experience. In 2023, more than 18 million IoT devices will be connected to the Internet [1] (as shown in Fig. 1).

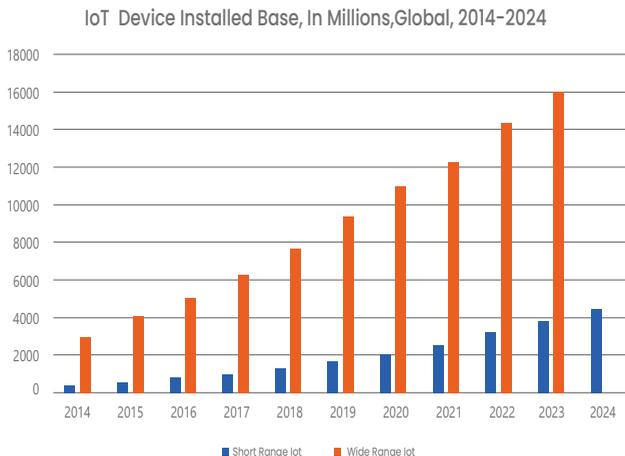


Figure 1. IoT devices installations [1]

At the same time, cybersecurity threats associated with IoT devices are on the increase. Top IoT threats are ranked according to their proportion rate (as shown in Fig. 2).

The occurrence of successful attacks targeting IoT systems can present dire consequences considering that these devices collect and store sensitive information. In addition to personal information, these devices contain sensors that take private images. From a legal and regulatory perspective, the need to ensure data privacy and confidentiality is paramount [3]. Machine learning, which encompasses machines, performs tasks without explicit programming, offers a promising way of detecting malware. Anomaly-based methods model the typical network behavior and identify abnormalities as plausible malware. This approach is advantageous as it offers a fairly effective method for dealing with both known and novel attacks [2]. The final hybrid approach

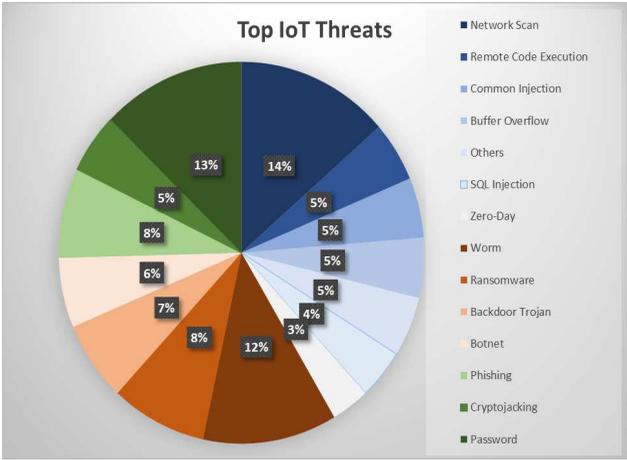


Figure 2. IoT Cybersecurity Threats [1]

combines both signature-based and anomaly-based methods. Machine learning techniques are routinely utilized to model the network behavior and identify anomalies. However, due to the high false positives, they often adopt the hybrid approach. By combining known malware signatures and network behavior, the ability of the intrusion detection system to detect malware accurately increases.

The application of deep learning in diverse big data fields has been successful due to the enhancement of computer processing power. Additionally, it offers an effective way of detecting attacks because of its high-level feature extraction capability. Essentially, deep learning uses a cascade of layers of processing units for extracting features from data, and each layer makes use of the output obtained from the prior layer as the input. Deep learning algorithms can adopt either supervised or unsupervised approaches. There seems to be an increasing adoption of deep learning methods as they can identify patterns and trends easily, can handle multi-variety and multi-dimensional data, and support continuous improvement. This paper intends to present an intelligent detection method for IoT systems that leverages effective classification, combined with various machine learning algorithms and a deep learning

algorithm, to detect benign and malicious activities.

The remainder of this paper is structured as follows: Section 2 presents the overview of related work. Dataset, Problem Statement and Proposed Framework is presented in Section 3. Subsections in Section 3 present more details of all the ML and DL algorithms that are used in this research. Section 4 presents the results of the proposed technique for intelligent detection of malicious activities in IoT networks. Section 5 details the conclusion of this research.

2 Related works

Different machine learning algorithms have been applied to anomaly detection in IoT systems and achieved positive outcomes [8], [9], [10], [11], [15], [16], [17], [19, 20, 21]. One of the most common approaches is the deep learning algorithm [8], [9].

Various datasets are available for examining the effectiveness of machine learning algorithms for detecting IoT-related intrusions [4], [5], [6], [22]. Koroniotis et al. developed one such dataset for the purposes of training and validating system credibility. The authors created a dataset called Bot-IoT, which integrated both simulated and legitimate internet of things traffic, including the different types of attacks. The study also encompassed presenting a test-bed environment for examining the current drawbacks of datasets, including capturing complete information about the network, accurate labeling, and emerging complex attacks. An evaluation of the BoT-IoT dataset using diverse machine learning and statistical methods compared with other datasets established adequate reliability [4].

Ullah and Mahmoud exploited a Botnet dataset from an existing one for detecting anomalous activity in IoT networks. The dataset had broader network and flow-based features, which were tested using diverse machine learning approaches, such as feature correlation, and recursive feature elimination. In addition to possessing the required accuracy levels, the dataset offers a good ground for analyzing anomalous activity detection models for IoT systems [5].

Sharafaldin et al. worked with an intrusion detection dataset. According to the researchers, most of the current datasets lack reliabil-

ity, especially in the face of emerging threats. After evaluating their dataset, the authors found that it exhibits reliability and accuracy when used together with machine learning algorithms to detect diverse attack categories [6]. However, this study did not focus on IoT-based systems, which is a key weakness as IoT networks tend to face unique attack threats. Nevertheless, the public availability of such datasets is imperative as it supports the creation and evaluation of IoT malicious detection models.

Shafiq et al. provided a specific detection method for malicious traffic in IoT systems. Despite the success of this study, the authors established that their approach did not detect some attacks accurately, especially keylogging. The performance of their approach for some machine learning algorithms such as decision trees and the random forest was also inadequate [1].

Dutta et al. followed the principle of stacked generalization to create an ensemble method that leverages deep learning models and a meta-classifier, which is the logistic regression [12]. The deep learning methods adopted encompassed the long short-term memory (LSTM) and the deep neural network (DNN) [12]. The approach utilized, encompassed two stages: the utilization of a Deep Sparse Auto-Encoder (DSAE) for feature engineering and a stacking ensemble for classification [12]. The evaluation of this approach showed that it is accurate in detecting network anomalies as compared to other state-of-the-art approaches. Abdullah et al. also illustrated the application of an ensemble method in the detection of network anomalies [13]. The system was based on dividing the input into different subsets based on the attack in question. Liu et al. presented a semi-supervised dynamic ensemble for detecting anomalies in IoT environments [14]. The algorithm combined mutual information criteria and semi-supervised extreme learning machine [14]. Experiments conducted on practical datasets showed that the proposed algorithm outperformed selected state-of-the-art approaches in terms of classification accuracy [14]. Evidently, ensemble approaches appear to be a promising direction of research mainly because of their ability to minimize biases and increase accuracy.

Tian et al. developed a distributed deep learning system for detecting web attacks on edge devices [8]. Based on the findings, the authors

established that deep learning is more effective in attack detection as compared to other approaches, especially when implemented in a distributed environment. Meidan et al. implemented deep encoders to detect IoT Botnet attacks [9]. The method entailed creating an algorithm that extracts behavior snapshots of the network and utilizes deep auto-encoders to detect abnormal traffic. The evaluation results showed that the method was effective in detecting attacks deployed using two widely known IoT-based Botnets. Thamilarasu and Chawla utilized a deep-learning algorithm to detect malicious traffic in IoT networks [11]. The system comprised network connection, anomaly detection, and mitigation modules. According to the authors, this system provided security as it served and facilitated interoperability between diverse network communication protocols utilized in IoT. The evaluation of the system demonstrated effectiveness and efficiency in detecting practical intrusions. Dutta et al. also illustrated the application of deep learning algorithms in the detection of anomalies in IoT systems [12]. Alhakami et al. implemented a non-parametric Bayesian approach to detect anomalies and identify intrusions [18]. This algorithm learned patterns of activities through a Bayesian-based MCMC inference for infinite bounded generalized Gaussian mixture models. The algorithm was tested and the outcomes showed that it was accurate in detecting diverse attacks.

More importantly, after detecting traffic anomalies, a system for classifying the attack is required. This classification is essential as it enables the implementation of the right controls. As a result, a more effective method is needed, which is the focus of this research. The idea is to adopt the ensemble approach of combining two or more machine learning methods to improve the accuracy for IoT attack detection.

3 Proposed Methodology

The application of machine learning algorithms to the detection of network anomalies and help in identifying and dealing with attacks is a recurrent topic in research literature [3]. Machine learning techniques developed often comprise two steps: training and testing [3]. Accordingly, the initial step encompasses identifying features or class attributes from

the training data. Afterwards, one has to identify a subset of attributed required to classify traffic either as normal or abnormal. This process is referred to as dimensionality reduction [3]. Once this is done, the model is trained using the training data and, thereafter, utilized in the classification of unknown data. During anomaly detection, the normal traffic pattern must be defined during the training process. Accordingly, when testing, the trained model is applied to new data, and every exemplar is classified as either anomalous or normal [3]. The same methodology has been followed to propose the ensemble approach in this paper and a step-by-step process is given below. According to Hasan et al., some of the attack and anomaly detection algorithms applied in IoT systems are logistic regression, decision trees, support vector machines, random forest, and artificial neural networks [7].

The growing research direction in developing machine learning systems for attack detection entails the utilization of ensemble algorithms [12, 13, 14]. The rationale behind ensemble learning is that a combination of two or more algorithms can produce better outcomes as compared to utilizing one algorithm alone. Stacking, bagging, and boosting are some of the common ensemble methods. Ensemble algorithms are ideal for classification as they minimize variance and biases hence boosting the accuracy of detection. In addition to whether it is supervised, semi-supervised, or unsupervised, the selection of an algorithm should also be based on its accuracy, recall, and precision.

3.1 Dataset to be used

The dataset used in the proposed study is known as UNSW-NB15 [4]. It is comprised of nine distinct types of attacks named as Fuzzers, Analysis, Backdoors, DoS, Exploits, Generic, Reconnaissance, Shellcode and Worms. It contains 11 attributes named as start time, last time, attack category, attack subcategory, protocol, source address, source port, a destination address, destination port, attack name, and attack reference. Training set of the data includes 175,341 records while the testing set contains 82,332 records. The dataset is also comprised of various network attributes such as protocol (proto), state, sent packets (spkts), discarded packets (dpkts) and attack category (attack cat), etc. The labeling of the dataset is done based on the attack cat-

egories (*attack_cat*). The target label attack category depicts either 0 or 1 depending on whether the record is normal or attacked.

3.2 Problem Statement

The problem is to identify which network, network nodes, protocols, services and states are prone to intruder attacks. It's a general classification problem where different feature variables are classified on the basis of attacks faced by the network, and upon accurate classification, future attacks can be prevented by making that particular feature node category more secure.

3.3 Proposed Framework

First of all, the training and testing sets are read into the workspace. To make the data ready for model implementation, it is passed through several preprocessing steps. Algorithm 1 given below describes all the steps involved in the proposed framework. The architecture of the proposed framework is shown in Fig. 3.

Algorithm 1.

Input: (*proto*), (*spkts*), (*dpkts*), (*attack_cat*) $\in F = FeatureSet$

Output: classification of normal or malicious intrusions

Initialisation:

step 1: training set and testing set are concatenated

step 2: missing values from the dataset are removed or replaced

step 3: categorical feature (*attack_cat*) selected

step 4: (*attack_cat*) is passed through "hot encoding" and "one hot encoding"

step 5: all data values are normalized

step 6: dataset is split into training and testing data

step 7: parameter tuninigs are done as given below:

Logistic Regression (Log. Reg.): *c* is a set equal to 0.1. *c* = array of *c* that maps to the best scores across every class. If *refit* is set to false, then for each class, the best *c* is the average of the *c*'s that correspond to the best scores for each fold. *C_* is of shape(*n_classes*) when the problem is binary. *Passive Aggressive Classifier (PAC):* *max* iterations are set to 50

K-Nearest Neighbor (KNN): *n_neighbors* are set to 3

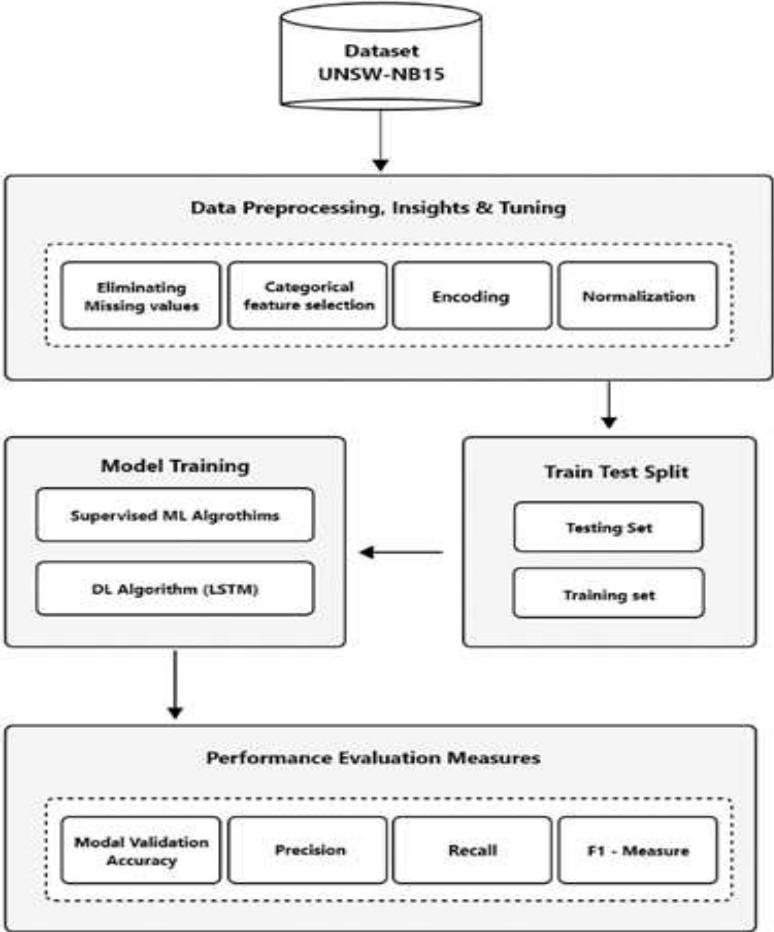


Figure 3. Architecture of Proposed Framework

Decision Tree Classifier (DT): random state is set to 1

Random Forest Classifier (Gini) (RF_Gini): n_estimators is set to 10

Random Forest (Information Gain) (RF_IFG): n_estimators is set to 10 and criterion = entropy

*Other algorithms are applied as default
end of algorithm*

3.4 Preprocessing and Data Preparation

Various preprocessing steps are applied so that the model can perform correctly on the dataset. Firstly, the missing values from the dataset are eliminated, then the dataset insights are taken in order to eliminate features that are not required or does not make any impact on the training of the model. Secondly, both the training and testing sets are concatenated for preprocessing process. Preprocessing can be done separately on training and testing partitions but it will have to be done twice which is why both the sets are concatenated. Finally, data normalization is performed to avoid overfitting.

3.5 Missing values

After the data have been concatenated, exploration is done to discover any missing values and then those values are replaced. The vast majority of the machine learning models that you want to employ will give you an error when you feed them *NaN* numbers. The best way is to avoid that and simply fill them with 0s.

3.6 Data insights and tuning

Although, the dataset is clean but it still needs further processing in terms of One-hot-encoding for categorical data. E.g., “service” consists of different types, we have ftp, http, and ‘-’ denoting not available or None. So we will need to treat it as a missing value as we will change it from ‘-’ to ‘None’ instead of dropping the whole column. Unnecessary features, e.g., “id” need to be removed as well.

3.7 Categorical features selection and encoding

The attack categories also known as attack cats in the dataset, are considered as a categorical feature in the proposed study because these attack categories contain all kinds of attacks that are experienced by

the network traffic. One hot encoding is applied on the feature attack category values so that these values can be inputted into the model as target values.

3.8 Data Normalization

The dataset normalization is done in order to avoid overfitting. Min-MaxScaler function is used to make sure that all the data values are between 0 and 1. As dataset size is large, so the MinMaxScaler function fits best in the scenario to normalize data. The normalization of the dataset will help machine learning algorithms to perform better or come to a conclusion faster. This is the reason, why it is a good practice to be followed before inputting such kind of data to the machine learning models.

3.9 Train and test splitting

The dataset is split into training and testing partitions in this case because our aim is to care for precision. For experimentation test data size is 30%. Parameter tunings are done and details about the use of various machine learning algorithms are given in Algorithm 1 given above. Some other machine learning algorithms such as Multi Nominal Naive Bayes (MNB), Gaussian Naive Bayes (GNB), Gradient Boosting (GB), Support Vector Machine (SVM) and deep learning algorithm (LSTM) are applied as default. The parameters' configuration of the ML and DL models is discussed below.

In the ML model MNB, alpha is set to 1.0, fit prior is set to true and class prior is set to none. Alpha is the parameter for additive smoothing of the data, and fit prior is used to whether learn or not learn about the prior class probabilities. We set class prior to none so that data prior class probabilities cannot be changed.

In the ML model GNB, the parameter prior is set to none. Variable smoothing parameter is also set so that zero probability can be handled. In the GB classifier, the parameters loss, learning rate equal to 0.1, and the number of estimators equal to 100 are set.

In the ML model SVM, the following parameters are set: parameter C equal to 1.0, kernel equal to rbf, gamma equal to scale, coefficient equal to 0.0, and shrinking equal to true. C is the parameter for regularization. The strength of the regularization goes down as C goes up.

The kernel is used to tell the algorithm what kind of kernel to use.

In the DL algorithm LSTM, batch size is set to 10, number of epochs are set to 100 with a validation split of 10 percent.

4 Results and Discussion

Finally, results are achieved through various machine learning models such as Logistic Regression, K-Nearest Neighbor (Lazy Algorithm), Decision Trees, Random Forest (Gini), Passive aggressive classifier, Multinomial naive bayes, Gaussian Naive Bayes, Gradient Boosting, SVM, Random Forest (Entropy or Information-gain), and a deep learning model such as LSTM.

The classification phase focuses on achieving best accuracy. Two types of accuracies are considered here, training accuracy and Cross-Validation (CV) accuracy. Table 1 shows training accuracies achieved by implementing the above-mentioned ML algorithm and LSTM that is a deep learning algorithm. The table also shows the execution time taken by each algorithm, validation accuracies, precision, recall, and F1 score.

It can be seen above that for our proposed method “Random Forest (Information Gain)” and “Passive Aggressive Classifier” provides the best cross-validation accuracies. These two machine learning algorithms have outperformed in the intelligent detection of benign and malicious attacks in the concerned IoT dataset used in our proposed detection method. Validation accuracies, precision, recall and F1 score for each algorithm used for classification of anomalies are explained in (equation 1, 2, 3, and 4) respectively.

K -fold cross validation accuracy is utilized to compute the accuracies of the models used in this study. 10 folds of the dataset are used to avoid overfitting.

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

where TP , TN , FP and FN stands for true positive, true negative, false positive, and false negative, respectively. TP is originally True as well as predicted True, TN is originally True, but predicted negatively

Table 1. Performance measures of the models

Model	Accuracy %	Execution Time	CV Accuracy	Precision	Recalls	F1-score
Logistic Regression	90.1	93.89	8.72	0.90	0.90	0.89
Decision Tree	99.74	93.89	91.1	0.99	0.99	0.99
Random Forest (Gini)	99.74	325.76	92.15	0.99	0.99	0.99
Random Forest (Information Gain)	99.74	388.32	92.2	0.99	0.99	0.99
Passive Aggressive Classifier	85.77	28.48	87.11	0.87	0.85	0.85
KNN	95.92	4040.04	89.71	0.95	0.95	0.95
Multi Nominal Naïve Bayes	74.45	15.24	74.33	0.80	0.74	0.74
Gaussian Naïve Bayes	50.5	12.12	50.46	0.79	0.50	0.44
Gradient Boosting	93.38	1234.75	91.57	0.93	0.93	0.93
SVM	90.28	367.53	88.78	0.90	0.90	0.90
LSTM	93 on 20 epochs	139.39	0.9359	0.9664	–	–

by the classifier. Figure 4 shows the k -fold accuracy of DL and ML models.

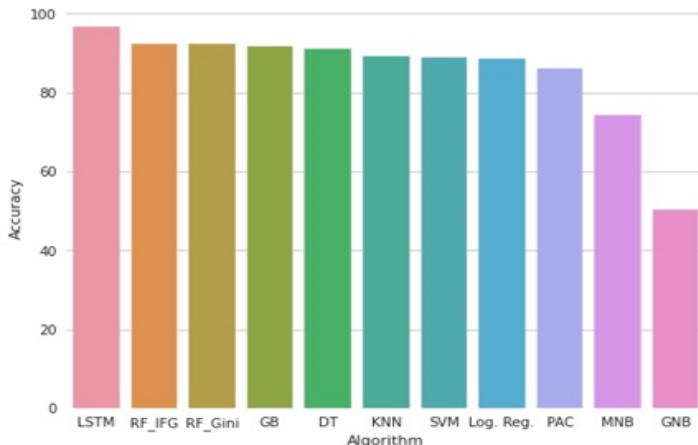


Figure 4. Validation Accuracies of various algorithms applied in the proposed study

Precision is percentage of correct predictions of a class among all predictions for that class.

$$Precision = \frac{[TP]}{[FP] + [TP]}, \quad (2)$$

where TP stands for true positive and FP stands for false positive. FP depicts the record which is originally false but predicted positive by the classifier. Figure 5 depicts precision score achieved by various machine and deep learning algorithms.

Recall is proportion of correct predictions of a class and the total number of occurrences of that class.

$$Recall = \frac{[TP]}{[FN] + [TP]}, \quad (3)$$

where TP stands for true positive and FN stands for false negative. FN depicts the record which is originally false and predicted false by the

classifier. In Fig. 6, precisions of various machine and deep learning algorithms are shown.

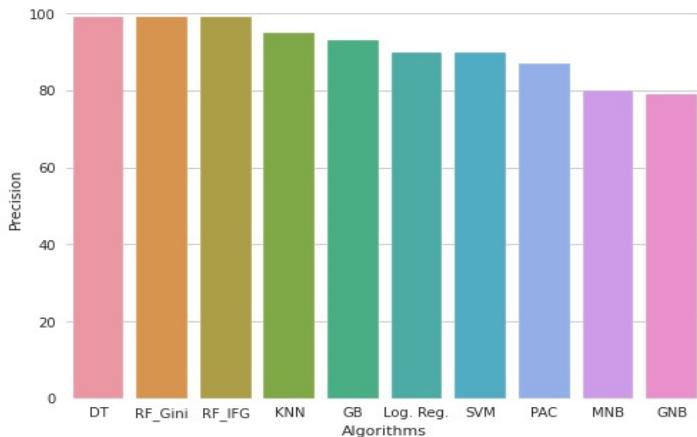


Figure 5. Precision of various algorithms applied in the proposed study

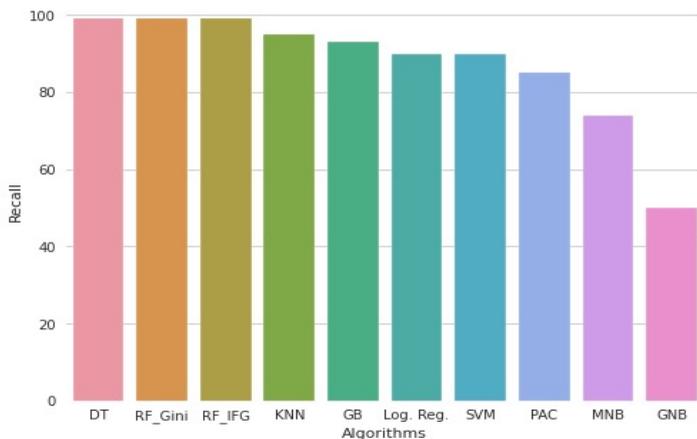


Figure 6. Recall of various algorithms applied in the proposed study

F1-score is a single metric combination of precision and recall.

$$F1 - score = 2 * \frac{[Recall] * [Precision]}{[Recall] + [Precision]} \tag{4}$$

Figure 7 depicts the F1-score achieved by machine and deep learning algorithms.

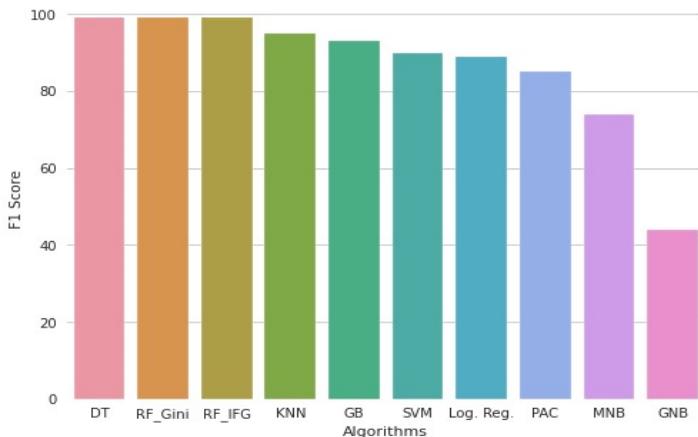


Figure 7. F1-score of various algorithms applied in the proposed study

For the evaluation, our proposed approach has compared with the state-of-art technique provided in [1] for training accuracy, execution time, Cross-Validation (CV) accuracy, precision, recall, and F1 score. Results are outstanding in the case of our proposed detection method and are shown in Fig. 8.

5 Conclusions

To classify benign and malicious intrusions in IoT devices is the big issue in this age of internet. In this research, Machine Learning algorithms and a Deep Learning method are being utilized for the intelligent detection of anomalies or malicious activities in IoT. Among many features being controlled during IoT network traffic, only some of them are responsible for the activation of malicious activities. This paper has

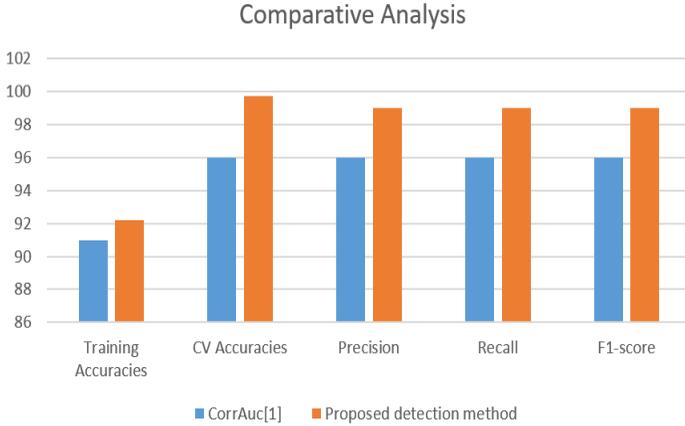


Figure 8. Comparative Analysis of CorrrAuc [1] and Proposed Approach

proposed an intelligent malicious detection technique that is basically an ensemble approach combined with various machine learning algorithms and deep learning algorithms. Apart from various algorithms being used in the experimentation, “Random Forest Information Gain” and “Passive aggressive classifier” provided the best cross-validation accuracies for anomaly detection. Moreover, after the evaluation, our proposed detection method has shown the best accuracies around 99%, which is more than any other state-of-art techniques developed in the literature.

Acknowledgement

The authors would like to thank Arab Open University, Saudi Arabia for supporting this study.

References

- [1] M. Shafiq, Z. Tian, A. K. Bashir, X. Du, and M. Guizani, “CorrAUC: A Malicious Bot-IoT Traffic Detection Method in IoT Network Using Machine-Learning Techniques,” *IEEE Internet of*

- Things Journal*, vol. 8, no. 5, pp. 3242–3254, March 1, 2021, DOI: 10.1109/JIOT.2020.3002255.
- [2] A. L. Buczak and E. Guven, “A survey of data mining and machine learning methods for cyber security intrusion detection,” *IEEE Communications Surveys and Tutorials*, vol. 18, no. 2, pp. 1153–1176, 2015.
 - [3] R. M. Alhajri, A. B. Faisal, and R. Zagrouba, “Survey for anomaly detection of IoT botnets using machine learning auto-encoders,” *Int J Appl Eng Res*, vol. 14, no. 10, pp. 2417–2421, 2019.
 - [4] N. Koroniotis et al., “Towards the development of realistic botnet dataset in the internet of things for network forensic analytics: Bot-Iot dataset,” *Future Generation Computer Systems*, vol. 100, pp. 779–796, 2019.
 - [5] I. Ullah and Q. H. Mahmoud, “A Technique for Generating a Botnet Dataset for Anomalous Activity Detection in IoT Networks,” in *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pp. 134–140, 2020.
 - [6] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, “Toward generating a new intrusion detection dataset and intrusion traffic characterization,” in *Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP 2018)*, pp. 108–116, 2018.
 - [7] M. Hasan et al., “Attack and anomaly detection in IoT sensors in IoT sites using machine learning approaches,” *Internet of Things*, vol. 7, Article ID: 100059, 2019. DOI: 10.1016/j.iot.2019.100059.
 - [8] Z. Tian, C. Luo, J. Qiu, X. Du, and M. Guizani, “A distributed deep learning system for web attack detection on edge devices,” *IEEE Transactions on Industrial Informatics*, vol. 16, no. 3, pp. 1963–1971, 2020.
 - [9] Y. Meidan, M. Bohadana, Y. Mathov, Y. Mirsky, A. Shabtai, D. Breitenbacher, and Y. Elovici, “A network-based detection of IoT botnet attacks using deep autoencoders,” *IEEE Pervasive Computing*, vol. 17, no. 3, pp. 12–22, 2018.
 - [10] V. Timcenko and S. Gajin, “Machine learning based network anomaly detection for IoT environments,” in *ICIST-2018 Conference*, 2018.

- [11] G. Thamilarasu and S. Chawla, “Towards deep-learning-driven intrusion detection for the internet of things,” *Sensors*, vol. 19, no. 9, pp. 1977, 2019. DOI: <https://doi.org/10.3390/s19091977>.
- [12] V. Dutta et al., “A deep learning ensemble for network anomaly and cyber-attack detection,” *Sensors*, vol. 20, no. 16, pp. 4583, 2020.
- [13] M. Abdullah et al., “Enhanced intrusion detection system using feature selection method and ensemble learning algorithms,” *International Journal of Computer Science and Information Security (IJCSIS)*, vol. 16, no. 2, pp. 48–55, 2018.
- [14] S. Liu et al., “A semi-supervised dynamic ensemble algorithm for IoT anomaly detection,” in *2020 International Conferences on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData) and IEEE Congress on Cybermatics (Cybermatics)*, pp. 264–269, 2020.
- [15] S. Egea, A. R. Manez, B. Carro, A. Sanchez-Esguevillas, and J. Lloret, “Intelligent IoT traffic classification using novel search strategy for fast based-correlation feature selection in industrial environments,” *IEEE Internet of Things Journal*, vol. 5, no. 3, pp. 1616–1624, 2018.
- [16] H. Zhang, G. Lu, M. T. Qassrawi, Y. Zhang, and X. Yu, “Feature selection for optimizing traffic classification,” *Computer Communications*, vol. 35, no. 12, pp. 1457–1471, 2012.
- [17] S. Su et al., “A correlation-change based feature selection method for IoT equipment anomaly detection,” *Applied Sciences*, vol. 9, no. 3, pp. 437, 2019.
- [18] W. Alhakami et al., “Network anomaly intrusion detection using a nonparametric bayesian approach and feature selection,” *IEEE Access*, vol. 7, pp. 52181–5219, 2019.
- [19] M. Shafiq et al., “Effective feature selection for 5G IM applications traffic classification,” *Mobile Information Systems*, vol. 2017, Article ID: 6805056, 12 pages, 2017.
- [20] I. H. Sarker et al., “Intrudtree: A machine learning based cyber security intrusion detection model,” *Symmetry*, vol. 12, no. 5, Article No. 754, 15 pages, 2020.

- [21] I. Ullah and Q. H. Mahmoud, “A two-level flow-based anomalous activity detection system for IoT networks,” *Electronics*, vol. 9, no. 3, Article No. 530, 2020. DOI: <https://doi.org/10.3390/electronics9030530>.
- [22] R. Ahmad and I. Alsmadi, “Machine learning approaches to IoT security: A systematic literature review,” *Internet of Things*, Article ID: 100365, 2021.

Saman Iftikhar, Danish Khan, Daniah Al-Madani,
Khattab M Ali Alheeti, Kiran Fatima

Received January 10, 2022
Revised July 24, 2022
Accepted September 1, 2022

Saman Iftikhar
Faculty of Computer Studies, Arab Open University, Saudi Arabia
E-mail: s.iftikhar@arabou.edu.sa

Danish Khan
Department of Computer Science, COMSATS University Islamabad
Wah Campus, Wah Cantt Pakistan
E-mail: danish56566@gmail.com

Daniah Al-Madani
Faculty of Computer Studies, Arab Open University, Saudi Arabia
E-mail: d.almadani@arabou.edu.sa

Khattab M Ali Alheeti
Computer Networking Systems Department,
College of Computer Sciences and Information Technology,
University of Anbar, Anbar, Iraq
E-mail: co.khattab.alheeti@uoanbar.edu.iq

Kiran Fatimah
TAFE, NSW Australia
E-mail: kiran.fatima4@tafensw.edu.au

Residual Neural Network in Genomics

Sara Sabba, Meroua Smara, Mehdi Benhacine, Loubna Terra,
Zine Eddine Terra

Abstract

Residual neural network (ResNet) is a Deep Learning model introduced by He et al. in 2015 to enhance traditional convolutional neural networks proposed to solve computer vision problems. It uses skip connections over some layer blocks to avoid vanishing gradient problem. Currently, many researches are focused to test and prove the efficiency of the ResNet on different domains such as genomics. In fact, the study of human genomes provides important information on the detection of diseases and their best treatments. Therefore, most of the scientists opted for bioinformatics solutions to get results in a reasonable time.

In this paper, our interest is to show the effectiveness of the ResNet model on genomics. For that, we propose two new ResNet models to enhance the results of two genomic problems previously resolved by CNN models. The obtained results are very promising and they proved the performance of our ResNet models compared to the CNN models.

Keywords: Deep Learning, genomics, convolutional neural network, companion, Residual neural network, super-enhancers, viral genomes.

1 Introduction

Machine Learning (ML) is one of the artificial intelligence fields, which is interested in the design and development of intelligent algorithms that learn and evolve with experiences to discover knowledge or make decisions (predictions) without being humanly guided or explicitly programmed to handle particular data. Indeed, the learning process begins with observations of data, experience, instructions, or examples to find the best model which can be able to make the best decisions in the future.

Deep Learning (DL) is the emerging technique of Machine Learning. Its basic concepts and models have derived from the Artificial Neural Network which mimics the activity of the nervous system of the human brain to intelligitize algorithms and avoid tedious human labor. DL has several computational models such as Deep fully connected neural networks (DNN), Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), Auto-encoder, Generative adversarial networks (GAN), Graph convolutional neural networks (GCN), Residual Neural Network (ResNet), etc. Most of them have provided their effectiveness in specific research areas such as computer vision [51, 52], natural language processing [53, 54, 55], and signal processing [56].

Currently, Deep Learning is an extremely active research area in bioinformatics [7, 15, 24, 26, 27, 37] due to the massive evolution of biological data. Its algorithms proved their efficiency in many critical life situations. They allow predicting many diseases, treatments, and biological phenomena from the analysis and interpretation of various types of data [1, 9, 10, 22, 24, 41]. In fact, most of the bioinformatics research is focused on Molecular biology which usually is called genomic. It is mainly interested in studying the cell at the molecular level, i.e., understanding the interactions between the different molecular systems of a cell, including the interactions between these macromolecules (DNA, RNA, and protein biosynthesis), as well as learning how these interactions are regulated.

In fact, many bioinformatics frameworks based on Deep Learning were developed in the literature to solve genomics problems. Xu et al. [37] proposed DeepEnhancer framework for predicting enhancers using convolutional neural networks (CNN). They used the FANTOM5 permissive enhancer dataset, JASPAR database and ENCODE cell type-specific enhancer datasets to train and test their model. Zhou et al. [40] developed a Deep Learning-based algorithmic framework, called DeepSEA to predict the noncoding-variant effects de novo from sequence. The proposed model is trained and tested on a regulatory sequence code from large-scale chromatin-profiling data. Alipanahi et al. [3] used also deep convolutional neural networks to develop DeepBind approach for predicting the sequence specificities of DNA- and RNA-

binding proteins. This approach is trained and tested on in vitro data, and it addressed many challenges we cite: (i) it can be applied to both microarray and sequencing data; (ii) it can tolerate a moderate degree of noise and mislabeled training data and (iii) it can train predictive models fully automatically, alleviating the need for careful and time-consuming hand-tuning. Likewise, SpliceFinder [34] and Splice2Deep [2] were designed to predict splice sites of human genomic using CNN model. The both works are trained and validated on some genomic sequences such as Homo sapiens, Oryza sativa japonica, Mus musculus, Drosophila melanogaster, and Danio rerio. In fact, there are so many critical frameworks worthy of our interest that we cannot cite them all.

Most of the genomic Deep Learning solutions are based on the CNN models. As is already known, CNNs are very useful in solving image classification and visual recognition problems. However, studies have shown that if the network has too many layers, we can observe the degradation of performance due to the vanishing or exploding gradient problem [57]. Accordingly, ResNet was introduced [42] to solve this problem by using skip connections (identity connections) or shortcuts (that create residual blocs) to jump over some layers which allow the network to retain what it has previously learned.

Recently, researchers are motivated to implement new genomic solutions using ResNet models we cite: Li et al. [44] developed ResPRE method to predict residue-level protein contacts using inverse covariance matrix of multiple sequence alignments. Sun et al. [45] proposed RNAcontact algorithm for predicting RNA inter-nucleotide 3D closeness. Shuvo et al. [46] introduced QDeep method to present new distance-based single-model quality estimation by harnessing the power of stacked deep ResNets. Zhang et al. [47] predicted Gene Expression from DNA Sequence. Zhang and Shen [49] proposed ThreaderAI to improve protein tertiary structure prediction. Kandel et al. [50] presented PURESNet model for predicting protein-ligand binding sites. Li and Xu [48] developed a new model of convolutional residual neural network for predicting protein structure using Inter-residue distance prediction. Wang et al. [43] proposed RPreS to predict RNA secondary structure profile. However, the number of these proposals remains modest compared to the CNN ones.

In this paper, we propose two new residual neural network models for two genomic problems. The first proposition aims for predicting super-enhancers on a genome scale, and the second aims for predicting viral genomes. Our purpose is to improve the results obtained by previous solutions based on CNN models [5, 59] and prove the effectiveness of ResNet models in genomic science. Moreover, there are three reasons behind this motivation: (i) first, ResNet was created to optimize the performance of CNN for avoiding the vanishing gradient problem, (ii) second, to the best of our knowledge, none of the literature research on super-enhancers or vital genome prediction is utilizing ResNet-based approach, and (iii) third, the obtained results proved the performance of our proposals compared to the CNN models.

2 Related works

2.1 Super-enhancers prediction

The prediction of super-enhancers (SEs) has prominent roles in biological and pathological processes. They play critical roles in the control of cell-type-specific genes programs, especially that related to the detection and progression of tumors [8, 14, 18, 32, 36]. SEs are defined as clusters of transcriptional enhancers. They are formed by binding of high levels of enhancer-associated chromatin features that drive high-level expression of genes encoding key regulators of cell identity [16, 26, 30].

The identification of SEs is based on the differences in their ability to bind markers of promoter transcriptional activity [32], including cofactors such as mediators (MED1, MED12) and cohesions (Nipbl, Smc1), histone modification markers (H3K27ac, H3K4me1, H3K4me3, H3K9me3), chromatin regulators (Brg1, Brd4, Chd7), and chromatin molecules (p300, CBP). Furthermore, Whyte et al. [35] indicated five embryonic stem cell (ESC) transcription factors to occupy super-enhancers (Oct4, Sox2, Nanog, Klf4, and Esrrb). However, there are many additional transcription factors that contribute to the control of ESCs [27, 29, 39]. In [14], authors tested ChIP-Seq data for fifteen additional transcription factors in ESCs and explored whether they occupy enhancers defined by Oct4, Sox2, and Nanog (OSN) co-occupancy. Their experiment results showed that six additional transcription fac-

tors (Nr5a2, Prdm14, Tcfcp2l1, Smad3, Stat3, and Tcf3) occupy both typical enhancers and super-enhancers and that all of them are enriched in super-enhancers.

Recently, many studies [23, 31, 32] proved that gene transcriptional dysregulation is one of the core tenets of cancer development that involves in noncoding regulatory elements, such as TFs, promoters, enhancers, SEs, and RNA polymerase II (Pol II). In particular, SEs play core roles in promoting oncogenic transcription to accelerate cancer development [4, 7, 32]. Recent research showed that cancer cells acquire super-enhancers at the oncogene, and cancerous phenotype relies on these abnormal transcription propelled by SEs [13, 25]. Accordingly, it is important to understand super-enhancers and their components since they control much disease-associated sequence variation that occurs in these regulatory elements [12, 14, 21] in large amounts of data in order to better understand biological processes. This knowledge can lead to discoveries that improve quality of life (i.e., designing more effective medical treatments or discovering certain severe illness in its early stages).

There are few bioinformatics works based on Machine Learning proposed to predict super-enhancers of the genomes. Authors of [19] implemented and compared six different Machine Learning models to identify key features of SEs and to investigate their relative contribution to the prediction. The six models include: Random Forest, Support Vector Machine, k-Nearest Neighbor, Adaptive Boosting, Naive Bayes, and Decision Tree. To validate their idea, they used 10-fold stratified cross-validation, independent datasets in four human cell-types and a set of publicly available data. Authors of [5] proposed a new computational method called DEEPSSEN for predicting super-enhancers based on a convolutional neural network. The proposed method is trained and tested on 36 SEs features, where 32 ones are used in [19], and 4 others are selected from ChIP-seq and DNase-seq datasets.

2.2 Viral genomes prediction

Viral metagenomics is the science that studies human, animal, and plant viral diseases. It consists of describing the total viral genome, or

virome for the discovery of new viruses. The results of metagenomics have allowed advances in diagnosis, molecular epidemiology, and viral evolution, and these studies have great relevance for re-evaluating concepts in pathology and, in particular, the biological role of viruses in an organism [60, 61, 62].

The detection of potential viral genomes in human biological samples is a crucial step in the viral metagenomics process. It currently represents an interesting problem in the field of bioinformatics. It aims to identify a human virome in DNA sequences extracted by a previous phase of metagenomics. Indeed, viruses are reservoirs and carriers of genes, suggesting that the human virome may have played a central role in human adaptation and evolution [64]. This importance reveals the need to update the methods used by this science.

In fact, there are some bioinformatics works based on Machine Learning proposed to identify viral DNA sequences. Ren et al. [64] implemented VirFinder based on the k-mers approach. Ren et al. [65] proposed a new computational method called DeepVirFinder based on a convolutional neural network. Tampuu et al. [59] enhanced the previous approach by proposing a parallel model called ViraMiner which is based on two CNN branches configured differently.

3 Proposed models

This section is divided into two parts. The first one presents the ResNet model proposed to predict super-enhancers, and the second part describes the ResNet model proposed to identify viral genomes.

3.1 ResSEN: Residual Neural Network for predicting super-enhancers

3.1.1 Datasets

The public database used to train and test our approach is used in [19] and [5] published previously. In fact, there are 36 features (see Table 1) incorporated in publicly available ChIP-seq and DNase-seq datasets of mouse embryonic stem cells (mESC) taken from Gene Expression Omnibus (GEO).

Table 1. Features of datasets used by [5] and our approach.

Super-enhancers data type	Features
Histone modifications	H3K27ac, H3K4me1, H3K4me3, H3K9me3
DNA hypersensitive site	DNaseI
RNA polymeraseII	Pol II
Transcriptional co-activating proteins	p300, CBP
P-TFEb subunit	Cdk9
Sub-units of Mediator complex	Med12, Cdk8
Chromatin regulators	Brg1, Brd4 and Chd7
Cohesin	Smc1, Nipbl
Subunits of Lsd1-NuRD complex	Lsd1, Mi2b
Histone deacetylase	H-DAC2, HDAC
Transcription factors	Oct4, Sox2, Nanog, Esrrb, Klf4, Tcfcp2l1, Prdm14, Nr5a2, Smad3, Stat3, Tcf3
Sequence signatures	AT content, GC content, phastCons, phastConsP, repeat fraction

The datasets contain 11100 samples. Among them, 1119 are positive and 9981 are negative. To train, test, and compare our ResSEN approach, we divided those samples into training datasets and test datasets, where 90% (i.e., 9990) are used for training and 10% (i.e., 1110) are used for performance testing (see Table 2).

Table 2. Division of samples.

Datasets	Samples size	Positive samples	Negative samples
Training datasets	9990	1006	8984
Test datasets	1110	113	997

Notice that the samples used in the validation phase are the same used in the test phase because the total number of samples is insufficient to be devised into three sub datasets.

3.1.2 ResSEN model

ResSEN model is composed of an input layer, a convolution layer, a pooling layer, two residual blocks and a fully connected layer.

a. Input layer

Thirty six (36) characteristics are used to predict the super-enhancers (see Table 1). So, there are 36 nodes in the input layer. The values of these nodes are normalized (using Eq. (1)) and standardized (using Eq. (2)) before they are transmitted to the next network layers.

$$y = (x - min) + (max - min), \quad (1)$$

where x is the input node value, and max, min are the maximum, minimum values between input nodes.

$$z = (y - mean) / standard_deviation, \quad (2)$$

where y is the normalized node value, $mean$ is calculated using Eq. (3), and $standard_deviation$ is calculated using Eq. (4):

$$mean = sum(y) / count(y), \quad (3)$$

$$standard_deviation = \sqrt{(sum((y - mean)^2)/count(x))}. \quad (4)$$

b. Convolutional layers

ResSEN is composed of five convolutional layers: i) a convolutional layer before the first residual block, and ii) two convolutional layers for each residual block ($2 \times 2 = 4$).

In the first convolutional layer we applied 64 filters of size 1×7 , followed by Max-pooling with pool-size 1×3 and stride 1. The first residual block has two convolutional layers, we applied 128 filters of size 1×3 in the first one, and 256 filters of the same size 1×3 in the second one. The second residual block has also two convolutional layers. In the first layer, we applied 256 filters of size 1×3 , while in the second layer, we applied 512 filters of the same size 1×3 .

Figure 3 illustrates the filters' parameters of the five convolutional layers.

Each convolutional layer is followed by a Batch Normalization (BN) layer (see Fig. 1) which is used to improve the speed, performance, and stability of deep neural networks [11], [17].

c. Activation layer

Deep learning usually employs a multilayer network and the gradient algorithm to train models, therefore it requires heavy computing, and the learning is often trapped into local minima. Currently, studies propose the rectified linear unit (ReLU) as the activation function to address this problem because its gradient is simple to compute, which allows the model to train easier, faster, and perform better [7].

Consequently, ResSEN uses ReLU as an activation function:

$$ReLU(x) = \max(0, x). \quad (5)$$

d. Add identity

For each residual block, ResSEN uses convolution block strategy to add the block's input to the block's output. This type of design requires that the block's output and its input have the same shape (size), so they can be added together. The output of the first block will be the

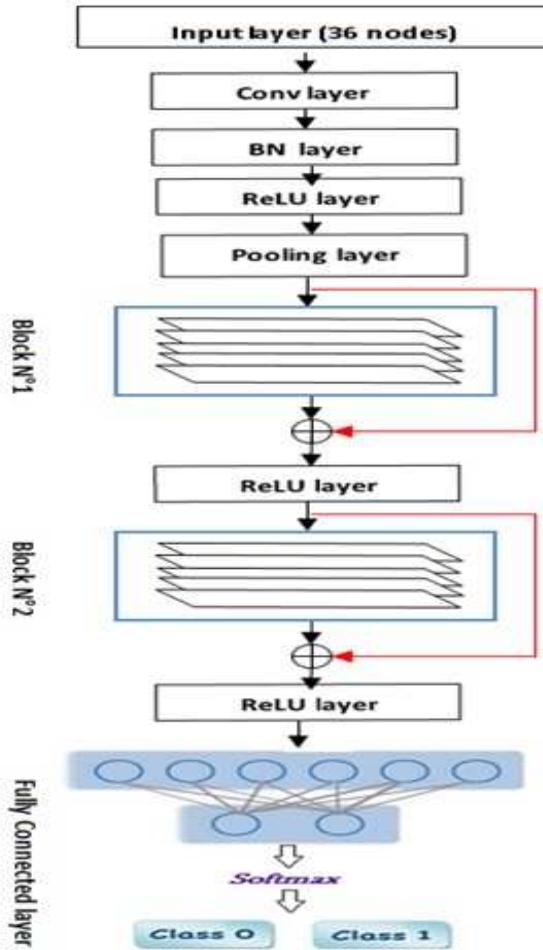


Figure 1. ResSEN model

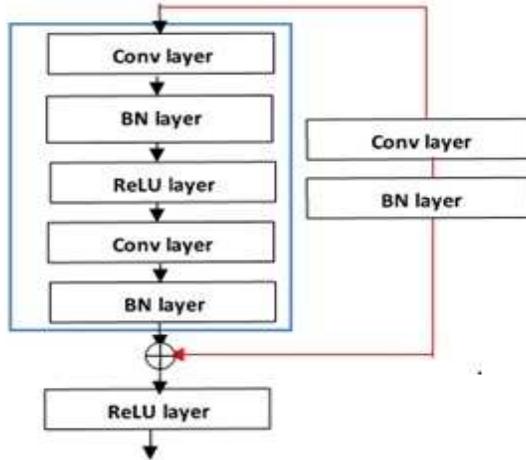


Figure 2. Residual Block of ResSEN

input of the second block and the output of the second block will be the input of the fully connected layer. The structure of each residual block is shown in Fig. 2.

To transform the block’s input into the desired shape, we introduced 256 convolutions (256 filters) of size 1×3 for the first residual block and 512 convolutions (512 filters) of size 1×3 for the second residual block (see Fig. 3).

e. Fully connected layer

The fully connected (FC) layer of the ResSEN is structured as follows:

- The number of input neurons is 17408.
- The activation function is ReLU.
- The number of output layer is 2 neurons.
- The function used to calculate the probability of the output classes is: Softmax (see Eq. (6)).

$$Softmax(x_j) = \max \frac{e^{x_i}}{\sum_j e^{x_i}}, j \in \{1, 2, \dots, k\}, \quad (6)$$

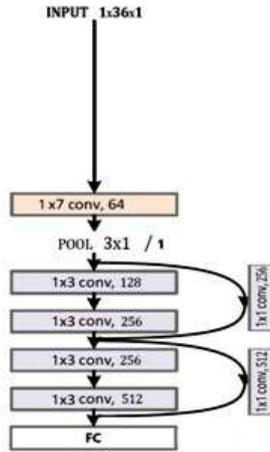


Figure 3. ResSEN convolutional layers parameters

where, k is the number of classes. Moreover, to obtain the predicted class A , we applied the *argmax* function to the *Softmax* function output:

$$A = \text{argmax}(\text{Softmax}(x_j)). \quad (7)$$

- So, if $A = 1$, the predicted class is positive, which means the presence of the super-enhancers in the genome;
- if $A = 0$, the predicted class is negative, which means the absence of super-amplifiers in the genome.

3.1.3 ResSEN training

ResSEN training is based on supervised learning, which consists of calculating the optimal weights using the input matrix D (the data samples) and the output matrix A (the desired outputs or the class labels) corresponding to D . D is a matrix of size $N \times 36$, and A is a binary matrix of size $N \times 1$, where N is the number of samples, which is set to 9900. $A[i] = 1$ if the corresponding sample represents

the super-enhancer class, otherwise, $A[i] = 0$. During the training phase, ResSEN uses the cross entropy loss function that measures the difference between the calculated output and the desired output (see Eq. (8)).

$$J(\Theta) = \frac{1}{n} \sum_{i=1}^n a^i \log(h_{\Theta}(x^i)) + (1 - a^i) \log(1 - h_{\Theta}(x^i)), \quad (8)$$

where, Θ is the set of parameters, n is the number of samples, a^i is the label of x^i , $h_{\Theta}(x^i)$ is the predicted label of x^i .

To update ResSEN weights, we used Backpropagation model and Adam method [20]. The latter is an adaptive learning rate optimization algorithm that is designed to improve the classical method of stochastic gradient descent (SGD) aiming at accelerating deep neural network learning. It automatically adapts the learning rate for each parameter by calculating adaptive estimates of moments.

3.2 ResVG: Residual Neural Network for predicting viral genomes

3.2.1 Datasets

The datasets used to train, validate, and test our approach are used in the work of Tampuu et al. [59] published previously. There is a set of metagenomic sequences taken from different samples such as skin, serum, and condyloma, obtained by merging and mixing 19 experiments. These last are divided into medium-sized ones (of 300 bp). The datasets contain 264049 samples. Among them, 5551 are positive (viral sequences) and 258498 are negative (non-viral sequences). To train, test, and compare our ResVG approach, we divided those samples into training, validation, and test datasets, where 80% (i.e., 211239) are used for training, 10% (i.e., 26405) are used for validation, and 10% (i.e., 26405) are used for performance testing (see Table 3).

3.2.2 ResVG model

ResVG model is composed of an input vector, convolutional layers, batch normalization layers, ReLU activation layers, a Max pooling

Table 3. Division of viral genome samples

Datasets	Samples size	Positive samples	Negative samples
Training datasets	211239	4466	206773
Test datasets	26405	551	25854
Validation datasets	26405	534	25871

layer, a residual block, a Global Average Pooling layer, and a fully connected layer.

a. Input data

ResVG uses DNA sequences of length 300 bp, each one is coded in binary on 4 bits. Indeed, there are 4 possible values (ACGT): A = 1000, C = 0100, G = 0010, and T = 0001. So, each input vector corresponds to a 1D sequence (an array) of length 300 with 4 channels (as shown in Fig. 4). This means that there are 1200 input neurons for the network.

b. Convolutional layers

ResVG is composed of three convolutional layers: i) a convolutional layer before the first residual block; ii) two convolutional layers for the residual block.

In the first convolutional layer, we applied 64 filters of size 1×7 , followed by Max-pooling with pool-size 1×4 and stride 1. The residual block has two convolutional layers; for both, we applied 1000 filters of size 1×11 .

As the first proposal, each convolutional layer is followed by a Batch Normalization layer (BN) and ReLU (rectified linear unit) activation layer.

c. Add identity

For a residual block, ResVG uses also convolution block strategy to add the block's input to the block's output. Therefore, to transform the ResVG block's input into the desired shape, we introduced 1000 convolutions (1000 filters) of size 1×11 (see Fig. 4).

d. Fully connected layer

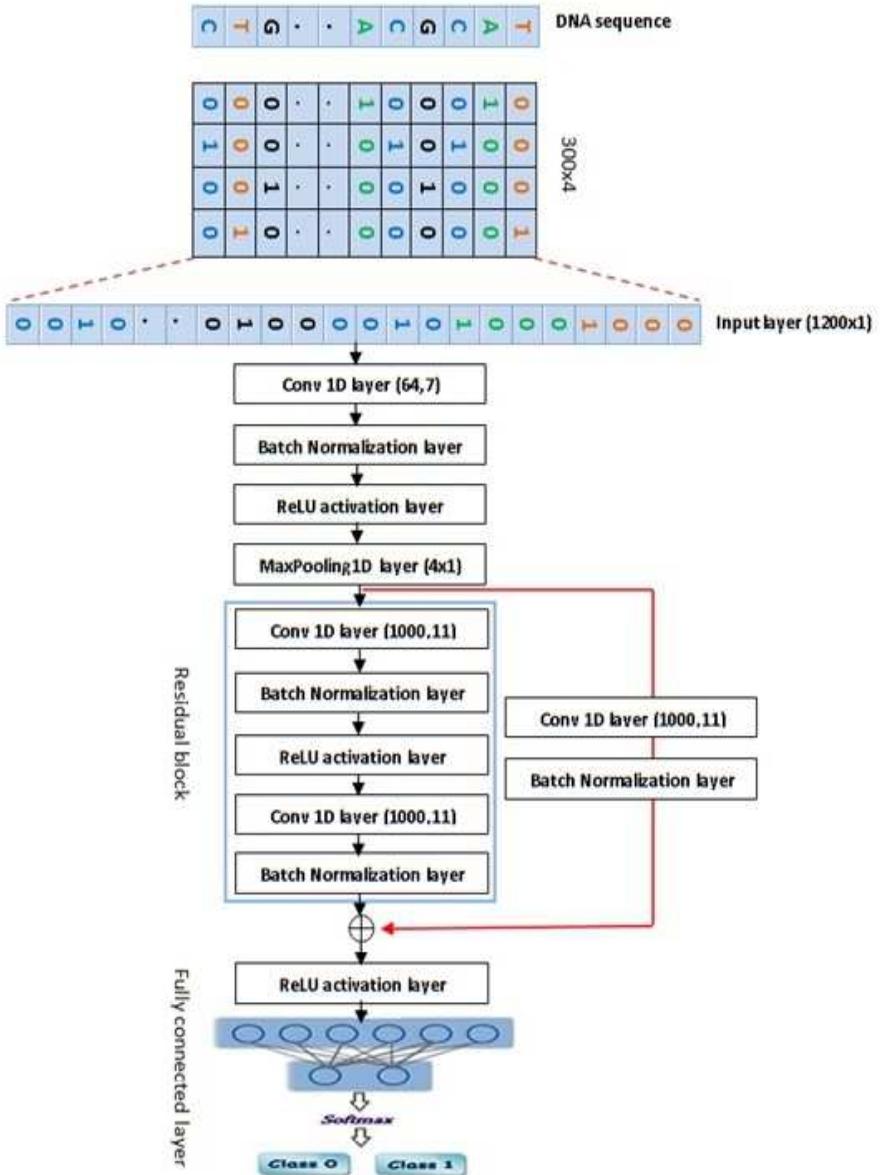


Figure 4. ResVG model

The fully connected (FC) layer of the ResVG is structured as follows:

- The number of input neurons is 1001.
- The activation function is ReLU.
- The number of output layer is 2 neurons.
- The function used to calculate the probability of the output classes is: Softmax (see Eq. (6)).

3.2.3 ResVG training

ResVG training is based on supervised learning, which consists of calculating the optimal weights using the input matrix D (the data samples) and the output matrix A (the desired outputs or the class labels) corresponding to D . D is a matrix of size $N \times 1200$ (300×4), and A is a binary matrix of size $N \times 1$, where N is the number of samples which is set to 211239. $A[i] = 1$ if the corresponding sample represents the viral class, otherwise, $A[i] = 0$.

During the training phase, ResVG uses the cross entropy loss function to measure the difference between the calculated output and the desired output, and Backpropagation model and Adam method [22] to update network weight's.

4 Experimental results and comparison

In the context of binary classification, the evaluation of models is based on some performance measures that are computed from the confusion matrix (see Table 4). Thus, to evaluate and compare our model's performance with those published by DeepSEN [5] and ViraMiner [59], we calculated accuracy, recall, precision, and F1-score for ResSEN model and accuracy, precision, and AUROC (TPR vs FPR) for ResVG model.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \quad (9)$$

$$TPR(TruePositiveRate)/Recall = \frac{TP}{TP + FN}, \quad (10)$$

$$Precision = \frac{TP}{TP + FN}, \quad (11)$$

Table 4. Confusion matrix

		Actual class	
		+	-
Predicted class	+	True Positives	True Negatives
	-	False Positives	False Negatives

$$F1Score = \frac{2 * Precision * Recall}{Precision + Recall}, \quad (12)$$

$$FPR(FalsePositiveRate) = \frac{FP}{TN + FP}. \quad (13)$$

Notice that the TPR and FPR are used by the AUROC curve to represent the separability degree between classes.

The best results obtained by testing the best models of ResSEN and ResVG are compared respectively with those of DeepSEN and ViramIner. Those comparisons are shown in Figs. 5 and 6.

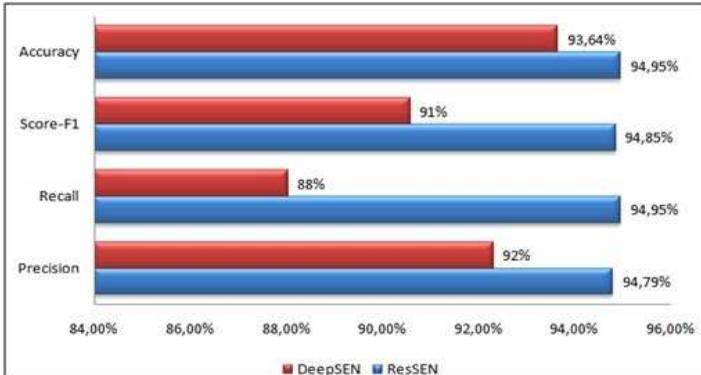


Figure 5. Performance comparison graph of ResSEN and DeepSEN in the validation and the test phases

In the DeepSEN paper, the authors proposed a model with three convolutional layers (followed each one by a pooling layer) and a fully connected layer. They mentioned that their model achieved an accuracy of 98% [5]. However, by checking the DeepSEN code published

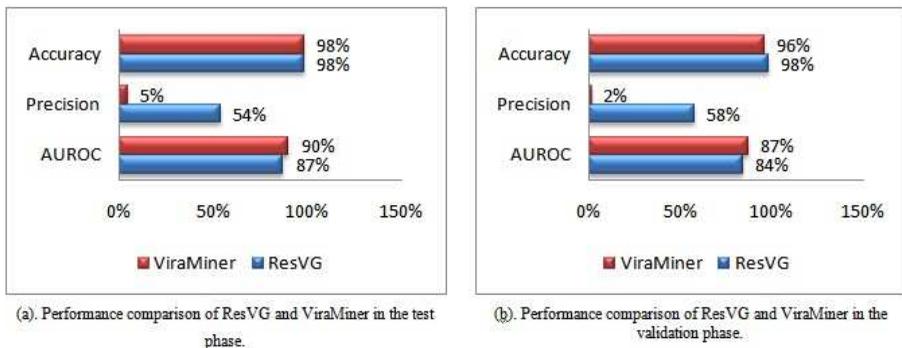


Figure 6. ResVG results

in [6], we found that they used 80% of samples in the training phase and all datasets (100% of samples) in the testing phase, which is wrong according to the learning strategy [33], [38]. Ideally, the model should be tested on samples that were not used in the training phase.

So, to ensure a fair comparison with our ResSEN model, we re-executed the DeepSEN using 90% of samples for training and 10% of samples for testing. The obtained results (see Fig. 8), show that the best model of the DeepSEN achieves an accuracy of 93,64% and a precision of 90%. However, in both cases (validation with all the datasets or with 10% of samples), we noticed the presence of the overfitting problem. The latter is clearly modeled in the accuracy and loss curves that we generated after re-executed DeepSEN (see Fig. 7). Knowing that, the blue and the orange curves represent the development of accuracy/loss in the training phase and in the validation phase, respectively.

Figure 8 shows the accuracy and loss curves of ResSEN model. In this case, there is no overfitting problem. We noticed a harmonization between the curves generated in training and test phases. Finally, the results shown in Figs. 5, 7, and 8 prove that our proposed model outperforms that of DeepSEN for the prediction of super-enhancers.

Furthermore, in ViraMiner paper, authors proposed a model with two branches. The first uses a single convolutional layer of 1000 filters followed by GlobalMaxPooling layer. The second branch also uses

a single convolutional layer of 1200 filters followed by GlobalAverage-Pooling layer. The results of the two branches are concatenated to find the inputs of the last FC layer.

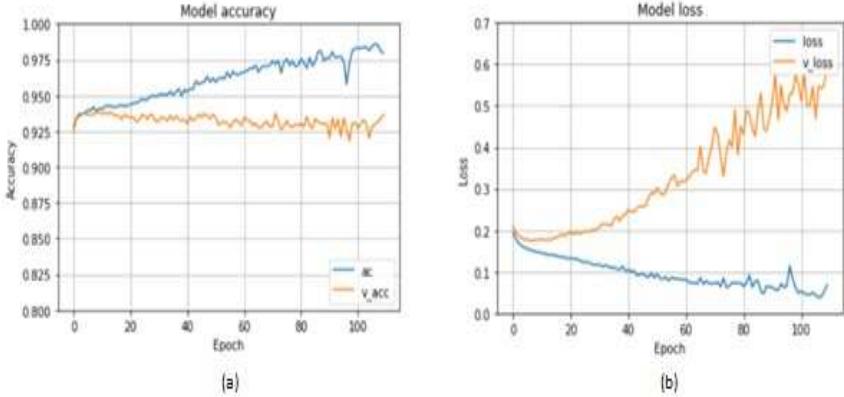


Figure 7. (a) DeepSEN accuracy curve, (b) DeepSEN loss curve

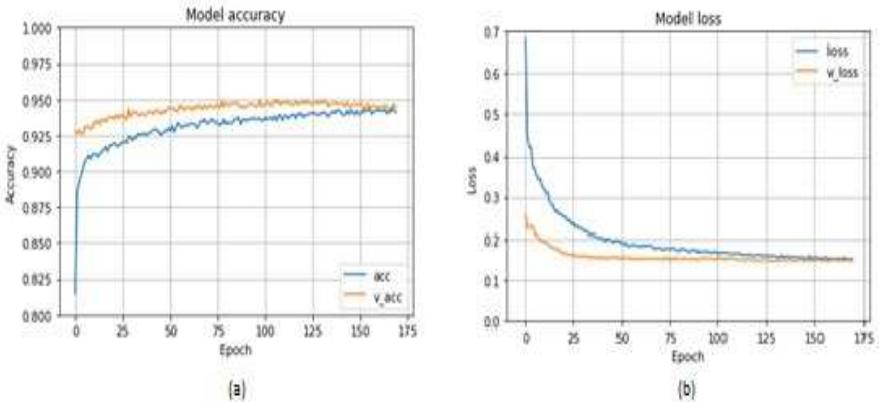


Figure 8. The feature vector set

To ensure a fair comparison with our ResVG model, we re-executed the ViraMiner model published in [66] using 80% of the samples for training, 10% for validation, and 10% for testing. The obtained results are shown in Figs. 6, 9, and 10.

Comparing the results of ViraMiner with our results, we noticed

that the AUROC curve of ViraMiner is the best; however, there is a big difference in the precision performance. Moreover, by analyzing the loss curves, we noticed that the ViraMiner model has a big overfitting problem compared to our Model. Finally, we can say that our ResVG model has optimized the prediction performance of viral genomes compared to the ViraMiner model, especially in the validation phase.

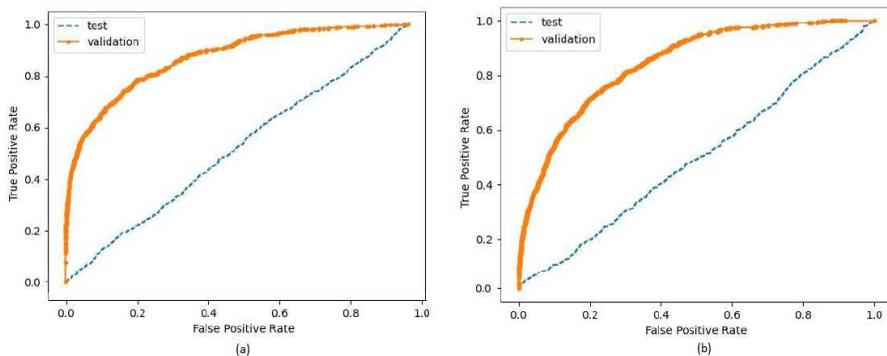


Figure 9. (a) The AUROC curve of ViraMiner, (b) The AUROC curve of ResVG

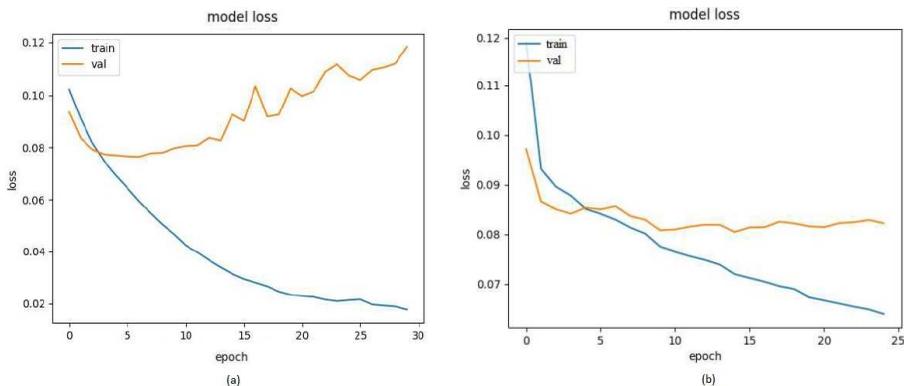


Figure 10. (a) The loss curve of ViraMiner, (b) The loss curve of ResVG

5 Conclusion

This paper is proposed to prove the performance of the ResNet model to solve genomic problems and to tackle the overfitting problem presented in the CNN models. Therefore, we proposed two ResNet models

called ResSEN and ResVG. The first model aims to predict the presence of super-enhancers on genome scale, it was tested and evaluated using 11100 samples composed each one of 36 features of mESC datasets taken from Gene Expression Omnibus (GEO). The second model aims to identify viral genomes, it was evaluated using 264049 metagenomic sequences of the size of 300 bp. The obtained results were compared respectively with those of two CNN models called DeepSEN [5] and ViraMiner [59] models. Comparisons showed that the overfitting problem is clearly disappeared in the ResSEN model and improved in the ResVG model. The final results showed also that the ResSEN is better than the DeepSEN for predicting super-enhancers and ResVG is better than the ViraMiner for identifying viral genomes but it can be more enhanced in the future by testing another optimized model of a CNN, like DenseNet or SENet. Finally, we conclude that the ResNet model can be a best solution for some genomic problems and it deserves to be tested in this domain.

References

- [1] M. Alazab et al., “COVID-19 Prediction and Detection Using Deep Learning,” *International Journal of Computer Information Systems and Industrial Management Applications*, vol. 12, pp. 168–181, 2020.
- [2] S. Albaradei et al., “Splice2Deep: An ensemble of deep convolutional neural networks for improved splice site prediction in genomic DNA,” *Gene X*, vol. 5, 2020.
- [3] B. Alipanahi, A. Delong, M. Weirauch, and B.J. Frey, “Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning,” *Nature Biotechnol*, vol. 33, pp. 831–838, 2015.
- [4] J.E. Bradner, D. Hnisz, and R.A. Young, “Transcriptional addiction in cancer,” *Cell*, vol. 168, pp. 629–643, 2017.
- [5] H. Bu, J. Hao, Y. Gan, et al., “DEEPSSEN: a convolutional neural network based method for super-enhancer prediction,” *BMC Bioinformatics*, vol. 20, 2019.
- [6] H. Bu, J. Hao, and Y. Gan, *DEEPSSEN code*, 2019. [Online]. Available: <https://github.com/1991Troy/DEEPSSEN>.

- [7] D. Chen, F. Hu, G. Nian, and T. Yang, “Deep Residual Learning for Nonlinear Regression,” *Entropy*, vol. 22, no. 2, Article ID: 193, 2020. <https://doi.org/10.3390/e22020193>.
- [8] S. Chen, Q. Jia, Q., Y. Tan, Y. Li, and F. Tang, “Oncogenic super-enhancer formation in tumorigenesis and its molecular mechanisms,” *Experimental & Molecular Medicine*, vol. 52, pp. 713–723, 2020.
- [9] T. Ching et al., “Opportunities and obstacles for deep learning in biology and medicine,” *Journal of the Royal Society Interface*, vol. 15, no. 141, Article ID: 20170387, 2018. DOI: 10.1098/rsif.2017.0387.
- [10] A. Esteva et al., “A guide to deep learning in healthcare,” *Nature Medecine*, vol. 25, no.1, 2019.
- [11] Y. Furusho and K. Ikeda, “ResNet and Batch-normalization Improve Data Separability,” in *Proceedings of Machine Learning Research*, vol. 101, pp. 94–108, 2019.
- [12] S.R. Grossman et al., “Identifying recent adaptations in large-scale genomic data,” *Cell*, vol. 152, pp. 703–713, 2013.
- [13] Y. He, W. Long, and Q. Liu, “Targeting Super-Enhancers as a Therapeutic Strategy for Cancer Treatment,” *Frontiers in Pharmacology*, vol. 10, 2019.
- [14] D. Hnisz, B.J. Abraham et al., “Super-enhancers in the control of cell identity and disease,” *Cell*, vol. 155, no. 4, pp. 934–947, 2013.
- [15] A. Holzinger and I. Jurisica, “Knowledge Discovery and Data Mining in Biomedical Informatics: The Future Is in Integrative, Interactive Machine Learning Solutions,” in *Interactive Knowledge Discovery and Data Mining in Biomedical Informatics* (Lecture Notes in Computer Science, vol. 8401), A. Holzinger and I. Jurisica, Eds. 2014.
- [16] J. Huang et al., “Dissecting super-enhancer hierarchy based on chromatin interactions,” *Nature Communications*, vol. 9, no. 943, 2018.
- [17] S. Ioffe and C. Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” *arXiv*, 2015. [Online]. Available: arXiv:1502.03167.

- [18] Y. Jia, W. Chng, and J. Zhou, “Super-enhancers: critical roles and therapeutic targets in hematologic malignancies,” *Journal of Hematology and Oncology*, vol.12, 2019.
- [19] A. Khan and X. Zhang, “Integrative modeling reveals key chromatin and sequence signatures predicting super-enhancers,” *Scientific Reports*, vol. 9, pp. 1–15, 2019. DOI: 10.1038/s41598-019-38979-9.
- [20] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations*, 2014.
- [21] T.I. Lee and R.A. Young, “Transcriptional regulation and its misregulation in disease,” *Cell*, vol. 152, pp. 1237–1251, 2013.
- [22] G. Litjens et al., “A survey on deep learning in medical image analysis,” *Medical Image Analysis*, vol. 42, pp. 60–88, 2017.
- [23] J. Lu et al., “MICAL2 mediates p53 ubiquitin degradation through oxidating p53 methionine 40 and 160 and promotes colorectal cancer malignance,” *Theranostics*, vol. 8, no. 19, pp. 5289–5306, 2018.
- [24] P. Mamoshina, A. Vieira, E. Putin, and A. Zhavoronkov, “Applications of deep learning in biomedicine,” *Molecular Pharmaceutics*, vol. 13, no. 5, pp. 1445–1454, 2016.
- [25] M. R. Mansour et al., “Oncogene regulation. An oncogenic super-enhancer formed through somatic mutation of a noncoding intergenic element,” *Science* (New York, N.Y.), vol. 346, no. 6215, pp. 1373–1377, 2014.
- [26] M.F et al., “Super-enhancers maintain renin-expressing cell identity and memory to preserve multi-system homeostasis,” *Journal Clinical Investigation*, vol. 128, no. 11, pp. 4787–4803, 2018.
- [27] S. Min, B. Lee, and S. Yoon, “Deep learning in bioinformatics,” *Brief Bioinform*, vol. 18, no. 5, pp. 851–869, 2017.
- [28] H.H. Ng and M.A. Surani, “The transcriptional and signalling networks of pluripotency,” *Nature cell biology*, vol. 13, pp. 490–496, 2011.
- [29] S.H. Orkin and K. Hochedlinger, “Chromatin connections to pluripotency and cellular reprogramming,” *Cell*, vol. 145, pp. 835–850, 2011.

- [30] J. Qu et al., “Functions and Clinical Significance of Super-Enhancers in Bone-Related Diseases,” *Frontiers in Cell and Developmental Biology*, vol. 8, 2020.
- [31] S. Sengupta and R.E. George, “Super-enhancer-driven transcriptional dependencies in cancer,” *Trends Cancer*, vol. 3, pp. 269–281, 2017.
- [32] F. Tang, Z. Yang, Y. Tan, and Y. Li, “Super-enhancer function and its application in cancer targeted therapy,” *NPJ Precision Oncology*, vol. 4, no. 2, 2020.
- [33] H. Wang and H. Zheng, “Model Validation, Machine Learning,” in *Encyclopedia of Systems Biology*, W. Dubitzky, O. Wolkenhauer, KH. Cho, H. Yokota, Eds. New York, NY: Springer, 2013. https://doi.org/10.1007/978-1-4419-9863-7_233.
- [34] R. Wang, Z. Wang, J. Wang, and S. Li, “SpliceFinder: ab initio prediction of splice sites using convolutional neural network,” *BMC Bioinformatics*, vol. 20, 2019.
- [35] W. Whyte et al., “Master transcription factors and mediator establish super-enhancers at key cell identity genes,” *Cell*, vol. 153, no. 2, pp. 307–319, 2013.
- [36] W. Xi, J.C. Murray, and Y. Jian, “Super-enhancers in transcriptional regulation and genome organization,” *Nucleic Acids Research*, vol.47, no. 22, pp. 11481–11496, 2019.
- [37] M. Xu, C. Ning, C. Ting, and J. Rui, “DeepEnhancer: Predicting enhancers by convolutional neural networks,” in *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, (Shenzhen), pp. 637–644, 2016.
- [38] Y. Xu and R. Goodacre, “On Splitting Training and Validation Set: A Comparative Study of Cross-Validation, Bootstrap and Systematic Sampling for Estimating the Generalization Performance of Supervised Learning,” *Journal of Analysis and Testing*, vol. 2, no. 3, pp. 249–262, 2018.
- [39] R.A. Young, “Control of the embryonic stem cell state,” *Cell*, vol. 144, pp. 940–954, 2011.

- [40] J. Zhou and O. Troyanskay, “Predicting effects of noncoding variants with deep learning–based sequence model,” *Nature Methods*, vol. 12, pp. 931–934, 2015.
- [41] H. Zilonga, T. Jinshanabc, W. Zimingb, Z. Kaiac, Z. Linga, and S. Qingling, “Deep learning for image-based cancer detection and diagnosis – A survey,” *Pattern Recognition*, vol. 83, pp. 134–149, 2018.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” *arXiv*, 2015. [Online]. Available: arXiv:1512.03385.
- [43] L. Wang, X. Zhong, S. Wang, S. et al., “A novel end-to-end method to predict RNA secondary structure profile based on bidirectional LSTM and residual neural network,” *BMC Bioinformatics*, vol. 22, 2021.
- [44] Y. Li., J. Hu, C. Zhang, D.J. Yu, and Y. Zhang, “ResPRE: high-accuracy protein contact prediction by coupling precision matrix with deep residual neural networks,” *Bioinformatics*, vol. 35, no. 22, pp. 4647–4655, Nov. 2019.
- [45] S. Sun, W. Wang, Z. Peng, and J. Yang, “RNA inter-nucleotide 3D closeness prediction by deep residual neural networks,” *Bioinformatics*, vol. 37, no. 8, pp. 1093–1098, Apr, 2021.
- [46] M.H. Shuvo, S. Bhattacharya, and D. Bhattacharya, “QDeep: distance-based protein model quality estimation by residue-level ensemble error classifications using stacked deep residual neural networks,” *Bioinformatics*, vol 36, pp. i285–i291, July, 2020.
- [47] Y. Zhang, X. Zhou, and X. Cai, “Predicting Gene Expression from DNA Sequence using Residual Neural Network,” *bioRxiv*, 2020. DOI: <https://doi.org/10.1101/2020.06.21.163956>.
- [48] J. Li and J. Xu, “Study of real-valued distance prediction for protein structure prediction with deep learning,” *Bioinformatics*, vol. 37, no. 19, pp. 3197–3203, Oct. 2021.
- [49] H. Zhang and Y. Shen, “Template-based prediction of protein structure with deep learning,” *BMC Genomics*, vol. 21, Supplement 11, Article number: 878, Dec. 2020. <https://doi.org/10.1186/s12864-020-07249-8>.

- [50] J. Kandel, H. Tayara, and K.T. Chong, “PUResNet: prediction of protein-ligand binding sites using deep residual neural network,” *J Cheminform*, vol. 13, no. 65, 2021.
- [51] S. Sharma, A. Juneja, and N. Sharma, “Using Deep Convolutional Neural Network in Computer Vision for Real-World Scene Classification,” in *IEEE 8th International Advance Computing Conference*, pp. 284–289, 2018.
- [52] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopadakis, “Deep Learning for Computer Vision: A Brief Review,” *Computational Intelligence and Neuroscience*, vol. 2018, Article ID 7068349, 13 pages, 2018. <https://doi.org/10.1155/2018/7068349>.
- [53] D. W. Otter, J. R. Medina, and J. K. Kalita, “A Survey of the Usages of Deep Learning for Natural Language Processing,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 2, pp. 604–624, Feb. 2021.
- [54] K.M. Tarwani and S. Edem, “Survey on recurrent neural network in natural language processing,” *Int. J. Eng. Trends Technol*, vol. 48, pp. 301–304, 2017.
- [55] W. Yin, K. Kann, M. Yu, and H. Schütze, “Comparative study of CNN and RNN for natural language processing,” *arXiv preprint*, arXiv:1702.01923, 2017.
- [56] H. Purwins et al., “Deep learning for audio signal processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 2, pp. 206–219, 2019.
- [57] C. Tian, Y. Xu, L. Fei, J. Wang, J. Wen, and N. Luo, “Enhanced CNN for image denoising,” *CAAI Transactions on Intelligence Technology*, vol. 4, no. 1, pp. 17–23, 2019.
- [58] L. Wen, Y. Dong, and L. Gao, “A new ensemble residual convolutional neural network for remaining useful life estimation,” *Math. Biosci. Eng*, vol. 16, no. 2, pp. 862–880, 2019.
- [59] A. Tampuu, Z. Bzhalava Z., J. Dillner, and R. Vicente, “ViraMiner: Deep learning on raw DNA sequences for identifying viral genomes in human samples,” *PLoS ONE*, vol. 14, no. 9. 2019.
- [60] S. Dávila-Ramos Sonia et al., “A Review on Viral Metagenomics in Extreme Environments,” *Frontiers in Microbiology*, vol. 10, 2019.

- [61] E.L. Delwart, “Viral metagenomics,” *Rev Med Virol.*, vol. 17, no. 2, pp. 115–131, 2007.
- [62] T.M. Santiago-Rodriguez and E.B. Hollister, “Potential Applications of Human Viral Metagenomics and Reference Materials: Considerations for Current and Future Viruses,” *Appl Environ Microbiol*, vol. 86, no. 22, e01794-20, Oct, 2020.
- [63] P. Bernardo, E. Albina, M. Eloit, and P. Roumagnac, “Pathology and viral metagenomics, a recent history,” *Medecine sciences*, vol. 29, no. 5, pp. 501–508, 2013.
- [64] J. Ren, N.A. Ahlgren, Y.Y. Lu, J.A. Fuhrman, and F. Sun, “VirFinder: a novel k-mer based tool for identifying viral sequences from assembled metagenomic data,” *Microbiome*, vol. 5, no. 69, 2017.
- [65] J. Ren et al., “Identifying viruses from metagenomic data by deep learning,” arXiv:1806.07810, 2020.
- [66] A. Tampuu, Z. Bzhalava, J. Dillner, and R. Vicente, *ViraMiner*, 2019. <https://github.com/NeuroCSUT/ViraMiner>.

Sara Sabba, Meroua Smara,
Mehdi Benhacine, Loubna Terra,
Zine Eddine Terra

Received November 28, 2021
Revised January 04, 2022
Accepted February 19, 2022

Sara Sabba

Department of Software Technologies and Information Systems,
Laboratory of Data Science and Artificial Intelligence(LISIA),
Abdelhamid Mahri University, Constantine 2, Algeria.
E-mail: sara.sabba@univ-constantine2.dz

Meroua Smara, Mehdi Benhacine, Loubna Terra, Zine Eddine Terra
Faculty of New Technologies of Information and Communication
Department of Software Technologies and Information Systems,
Abdelhamid Mahri University, Constantine 2, Algeria.
E-mail: meroua.smara@univ-constantine2.dz
E-mail: mehdi.benhacine@univ-constantine2.dz
E-mail: loubna.terra@univ-constantine2.dz
E-mail: zineeddine.terra@univ-constantine2.dz

Privacy and Reader-first Authentication in Vaudenay's RFID Model with Temporary State Disclosure

Ferucio Laurențiu Țiplea, Cristian Hristea, Rodica Bulai

Abstract

Privacy and mutual authentication under corruption with temporary state disclosure are two significant requirements for real-life applications of RFID schemes. This paper proposes two practical RFID schemes that meet these requirements. They differ from other similar schemes in that they provide reader-first authentication. Regarding privacy, our first scheme achieves destructive privacy, while the second one – narrow destructive privacy in Vaudenay's model with temporary state disclosure. To achieve these privacy levels, we use Physically Unclonable Functions (PUFs) to assure that the internal secret of the tag remains hidden from an adversary with invasive capabilities. Both of our schemes avoid the use of random generators on tags. Detailed security and privacy proofs are provided.

Keywords: Computer security, authentication, privacy, cryptography, PUF, RFID system.

MSC 2010: 94A60, 94A62.

1 Introduction

Radio Frequency Identification (RFID) refers to a technology whereby digital data encoded in *RFID tags* is transmitted to a *reader* via radio waves. A *back-end system*, which has an online database, is securely connected to the reader to collect, filter, process, and manage RFID data. It also stores complete information associated with the RFID tags in order to be able to authenticate them.

The RFID technology has been implemented in many significant areas such as toll collection systems, identification and tracking of various kinds of objects, consumer products, or access control. With the increasing usages to healthcare, electronic passports, and personal ID cards, the potential security threats and compliance risks have become enormous. In such a context, the need for secure and private communication protocols between reader and tags becomes crucial. Moreover, when developing such protocols, account must be taken of the adversary model to which they should resist. A widely accepted adversary model was proposed in [1],[2], now called *Vaudenay's RFID model*. According to it, an adversary can create legitimate or illegitimate tags, draw one or more tags according to some chosen probability distribution, release drawn tags, launch protocol instances with drawn tags, send messages to reader or drawn tags, corrupt drawn tags to retrieve their internal states, or get the result of a completed protocol instance.

Vaudenay's model classifies adversaries into eight classes and provides, consequently, eight levels of RFID privacy. Among these, destructive privacy (corruption destroys the tag) together with reader-first authentication under corruption with temporary state disclosure plays an important role in practice. For instance, tag destruction under corruption is an important requirement when the tag is used for access control. Likewise, the disclosure of temporary state under tag corruption is a serious threat in practice. Reader-first authentication [3] assures that the tag will give its private data only when it authenticates the reader. Therefore, tag tracking and data theft are prevented when the reader is fake. All these together mean that we need RFID schemes that provide destructive privacy and reader-first authentication under corruption with temporary state disclosure.

Contribution. When Vaudenay's model was proposed, it was not very clear whether the tag corruption reveals the permanent state or the full (permanent and temporary) state of the tag. Later, this aspect was clarified and it was shown that the mutual authentication protocols proposed in [2] do not achieve the claimed privacy level under corruption with temporary state disclosure. Additionally, this does not even happen [4] with newer protocols like those in [5],[6].

In this paper, we provide two mutual authentication RFID schemes

that achieve destructive and narrow destructive privacy in Vaudenay’s model with temporary state disclosure. Moreover, in our schemes, the tag authenticates first the reader (this is the reader-first approach [3]), which guarantees the information goes from tag to a trusted reader.

It is known that no privacy level can be achieved with ordinary tags (tags that only run cryptographic primitives) under mutual authentication and corruption with temporary state disclosure [7]. Therefore, the schemes we propose are based on *PUF tags*, that is tags endowed with *physically unclonable functions* (PUFs), a novel class of hardware security primitives that are in use for a while. The security proofs we provide to our schemes are very detailed. We also elaborate on the tag-reader desynchronization problem.

Related work. The pseudo-random function (PRF) based RFID scheme in [2] achieves weak privacy and mutual authentication in Vaudenay’s model. It is straightforward to see that the proof in [2] works even in the case of corruption with temporary state disclosure. The first PUF-based RFID scheme that achieves destructive privacy and mutual authentication in Vaudenay’s model (where corruption does not disclose the temporary state of tags) was proposed in [4], as an extension of the scheme in [8], [9] (that only achieves unilateral authentication).

In [5], [6], two PUF-based RFID schemes have been proposed and claimed that they achieve (narrow) destructive privacy and mutual authentication in Vaudenay’s model with temporary state disclosure. Unfortunately, neither of them reaches even narrow forward privacy [4].

Paper structure. The paper consists of six sections, the first one being the introduction. The basic terminology and notation used throughout this paper is introduced in Sections 2 and 3. Our first RFID scheme, that achieves destructive privacy and mutual authentication in Vaudenay’s model with temporary state disclosure, is presented in Section 4. In the fifth section, we introduce our second RFID scheme that achieves narrow destructive privacy and mutual authentication in the same model. The last section concludes the paper.

2 Basic definitions and notation

We fix here the basic terminology and notation for our paper.

Probabilistic polynomial time algorithms and negligible functions. We use *probabilistic polynomial time* (PPT) algorithms \mathcal{A} as defined in [10]. If \mathcal{O} is an oracle, then $\mathcal{A}^{\mathcal{O}}$ denotes that \mathcal{A} has oracle access to \mathcal{O} . When the oracle \mathcal{O} implements some function f , we simply write \mathcal{A}^f to denote that \mathcal{A} has oracle access to f . This means that whenever \mathcal{A} sends a value x to the oracle, it gets back $f(x)$.

If A is a set, then $a \leftarrow A$ means that a is uniformly at random chosen from A . If \mathcal{A} is a probabilistic algorithm, then $a \leftarrow \mathcal{A}$ means that a is an output of \mathcal{A} for some given input.

The asymptotic approach to security makes use of security parameters, denoted by λ in our paper. A positive function $f(\lambda)$ is called *negligible* if, for any positive polynomial $poly(\lambda)$, there exists n_0 such that $f(\lambda) < 1/poly(\lambda)$, for any $\lambda \geq n_0$.

Pseudo-random functions. Let ℓ_1 and ℓ_2 be two polynomials with positive values. Given a set \mathcal{K} of *keys* and $\lambda \in \mathbb{N}$, define $\mathcal{K}_\lambda = \{K \in \mathcal{K} \mid |K| = \lambda\}$. A *family of functions* indexed by \mathcal{K} is a construction $F = (F_K)_{K \in \mathcal{K}}$, where F_K is a function from $\{0, 1\}^{\ell_1(|K|)}$ to $\{0, 1\}^{\ell_2(|K|)}$. We also define $U_\lambda = \{f \mid f : \{0, 1\}^{\ell_1(\lambda)} \rightarrow \{0, 1\}^{\ell_2(\lambda)}\}$ and $U = (U_\lambda)_\lambda$.

We say that F is *computationally indistinguishable* from U if, for any PPT algorithm \mathcal{A} with oracle access to functions, its *advantage* $Adv_{\mathcal{A}, F}^{prf}(\lambda) = |P(1 \leftarrow \mathcal{A}^{F_K}(1^\lambda) : K \leftarrow \mathcal{K}_\lambda) - P(1 \leftarrow \mathcal{A}^g(1^\lambda) : g \leftarrow U_\lambda)|$ is negligible (as a function of λ).

$F = (F_K)_{K \in \mathcal{K}}$ is called a *pseudo-random function* (PRF) if it is efficiently computable and computationally indistinguishable from U .

Physically unclonable functions. A *physically unclonable function* (PUF) can be seen as a physical object that, when queried with a challenge x , generates a response y that depends on both x and the specific physical properties of the object. PUFs are typically assumed to be *physically unclonable* (it is infeasible to produce two PUFs that cannot be distinguished based on their challenge/response behavior), *unpredictable* (it is infeasible to predict the response to an unknown challenge), and *tamper-evident* (any attempt to physically access the PUF irreversibly changes the challenge/response behavior).

From a theoretical point of view, a PUF (sometimes called *ideal PUF*) is a physical object with a challenge/response behavior that implements a function $P : \{0, 1\}^p \rightarrow \{0, 1\}^k$, where p and k are of polyno-

mial size in λ , such that P is computationally indistinguishable from U , and any attempt to physically tamper with the object implementing P results in the destruction of P (P cannot be evaluated any more).

3 RFID schemes

From an informal point of view, an RFID system [11],[12] consists of a *reader*, a set of *tags*, and a *communication protocol* between reader and tags. The reader is a transceiver that has associated a database that stores information about tags. Its task is to identify *legitimate tags* (that is, tags with information stored in its database) and to reject all the other incoming communication. The reader and its database are trusted entities, and the communication between them is secure. A tag is a transponder device with much more limited computation capabilities than the reader. Depending on tag, it can perform simple logic operations, symmetric key, or even public key cryptography. Each tag has a *permanent* (or *internal*) *memory* that stores the state values, and a *temporary* (or *volatile*) *memory* that can be viewed as a set of *volatile variables* used to carry out the necessary computations.

RFID schemes. Let \mathcal{R} be a *reader identifier* and \mathcal{T} be a set of *tag identifiers* whose cardinal is polynomial in some security parameter λ . An *RFID scheme over* $(\mathcal{R}, \mathcal{T})$ [1],[2] is a triple $\mathcal{S} = (SetupR, SetupT, Ident)$ of PPT algorithms, where:

1. $SetupR(\lambda)$ inputs a security parameter λ and outputs a triple (pk, sk, DB) consisting of a key pair (pk, sk) and an empty database DB . pk is public, while sk is kept secret by reader;
2. $SetupT(pk, ID)$ initializes the tag identified by ID . It outputs an initial tag state S and a secret key K . A triple $(ID, f(S), K)$ is stored in the reader's database DB , where f is a public function that extracts some information from tag's initial state S ;
3. $Ident(pk; \mathcal{R}(sk, DB); ID(S))$ is an interactive protocol between the reader identified by \mathcal{R} (with its private key sk and database DB) and a tag identified by ID (with its state S) in which the reader ends with an output consisting of ID or \perp . The tag may end with no output (*unilateral authentication*), or it may end with an output consisting of OK or \perp (*mutual authentication*).

$SetupR(\lambda)$ “creates” a reader \mathcal{R} and initializes it, $SetupT(pk, ID)$ “creates” a tag \mathcal{T}_{ID} , initializes it with an initial tag state, and also registers this tag with the reader by storing some information about it in the reader’s database.

The *correctness* of an RFID scheme means that, regardless of how the system is set up, after each complete execution of the interactive protocol between the reader and a legitimate tag, the reader outputs tag’s identity with overwhelming probability. For mutual authentication of RFID schemes, *correctness* means that the reader outputs the tag’s identity, and the tag outputs *OK* with overwhelming probability.

An *RFID system* is an instantiation of an RFID scheme.

Adversaries. The two most basic security requirements for RFID schemes are *authentication* and *untraceability*. To formalize them, the concept of an *adversary model* is needed. There have been several proposals for this, such as [1],[2],[13]–[18]. One of the most influential, which we follow in this paper, is *Vaudenay’s model* [1],[2]. We recall below this model as in [4]. Thus, we assume first that some oracles the adversary may query share and manage a common list of tags $ListTags$, which is initially empty. This list includes exactly one entry for each tag created and active in the system. A tag entry consists of several fields with information about the tag, such as: the (permanent) identity of the tag (which is an element from \mathcal{T}), the temporary identity of the tag (this field may be empty saying that the tag is *free*), a bit value saying whether the tag is legitimate (the bit is one) or illegitimate (the bit is zero). When the temporary identity field is non-empty, its value uniquely identifies the tag, which is called *drawn* in this case. The adversary may only interact with drawn tags by means of their temporary identities.

The oracles an adversary may query are:

1. $CreateTag^b(ID)$: Creates a free tag \mathcal{T}_{ID} with the identifier ID by calling the algorithm $SetupT(pk, ID)$ to generate a pair (K, S) . If $b = 1$, $(ID, f(S), K)$ is added to DB , and the tag is considered *legitimate*; otherwise ($b = 0$), the tag is considered *illegitimate*. Moreover, a corresponding entry is added to $ListTags$;
2. $DrawTag(\delta)$: This oracle chooses a number of free tags according to the distribution δ , let us say n , and draws them. That is, n

- temporary identities $vtag_1, \dots, vtag_n$ are generated, and the corresponding tag entries in $ListTags$ are filled with them. The oracle outputs $(vtag_1, b_1, \dots, vtag_n, b_n)$, where b_i specifies whether the tag $vtag_i$ is legitimate or not;
3. $Free(vtag)$: Removes the temporary identity $vtag$ in the corresponding entry in $ListTags$, and the tag becomes free. The identifier $vtag$ will no longer be used. We assume that when a tag is freed, its temporary state is erased;
 4. $Launch()$: Launches a new protocol instance and assigns a unique identifier to it. The oracle outputs the identifier;
 5. $SendReader(m, \pi)$: Outputs the reader's answer when the message m is sent to it as part of the protocol instance π . When m is the empty message, abusively but suggestively denoted by \emptyset , this oracle outputs the first message of the protocol instance π , assuming that the reader does the first step in the protocol;
 6. $SendTag(m, vtag)$: outputs the tag's answer when the message m is sent to the tag referred to by $vtag$. When m is the empty message, this oracle outputs the first message of the protocol instance π , assuming that the tag does the first step in the protocol;
 7. $Result(\pi)$: Outputs \perp if in session π the reader has not yet made a decision on tag authentication (this also includes the case when the session π does not exist), 1 if in session π the reader authenticated the tag, and 0 otherwise (this oracle is both for unilateral and mutual authentication);
 8. $Corrupt(vtag)$: Outputs the current permanent (internal) state of the tag referred to by $vtag$, when the tag is not involved in any computation of any protocol step (that is, the permanent state before or after a protocol step).

We emphasize that $Corrupt$ does not return snapshots of the tag's memory during its computations. When the $Corrupt$ oracle returns the full state, we will refer to this model as being *Vaudenay's model with temporary state disclosure*.

Now, the adversaries are classified into the following classes, according to the access they get to these oracles:

- *Weak adversaries*: they do not have access to the $Corrupt$ oracle;

- *Forward adversaries*: once they access the *Corrupt* oracle, they can only access the *Corrupt* oracle;
- *Destructive adversaries*: after querying $Corrupt(vtag)$ and obtaining the corresponding information, the tag identified by $vtag$ is destroyed (marked as destroyed in $ListTags$), and the temporary identifier $vtag$ will no longer be available. The database DB will still keep the record associated to this tag (the reader does not know the tag was destroyed). As a consequence, a new tag with the same identifier cannot be created;
- *Strong adversaries*: there are no restrictions on the use of oracles.

Orthogonal to these classes, there is the class of *narrow* adversaries that do not have access to the *Result* oracle. We may now combine the narrow constraint with any of the previous constraints in order to get another four classes of adversaries, *narrow weak*, *narrow forward*, *narrow destructive*, and *narrow strong*.

Security. Now we are ready to introduce the *tag* and *reader authentication* properties as proposed in [1], [2], simply called the *security* of RFID schemes. First of all, we say that a tag \mathcal{T}_{ID} and a protocol session π had a *matching conversation* if they exchanged well interleaved and faithfully (but maybe with some time delay) messages according to the protocol, starting with the first protocol message but not necessarily completing the protocol session. If the matching conversation leads to tag authentication, then it will be called a *tag authentication matching conversation*; if it leads to reader authentication, it will be called a *reader authentication matching conversation*.

Now, the tag authentication property is defined by means of an experiment that a challenger sets up for a *strong* adversary \mathcal{A} (after the security parameter λ is fixed). In the experiment, the adversary is given the public parameters of the scheme and is allowed to query the oracles. If there has been a session in which the reader has authenticated an uncorrupted tag without a tag authentication matching conversation, then the experiment returns 1 (or 0 otherwise).

The advantage of \mathcal{A} in the experiment $RFID_{\mathcal{A},S}^{t_auth}(\lambda)$ is defined as

$$Adv_{\mathcal{A},S}^{t_auth}(\lambda) = Pr(RFID_{\mathcal{A},S}^{t_auth}(\lambda) = 1).$$

An RFID scheme \mathcal{S} achieves *tag authentication* if $Adv_{\mathcal{A},\mathcal{S}}^{t,auth}$ is negligible, for any strong adversary \mathcal{A} .

The experiment for reader authentication, denoted $RFID_{\mathcal{A},\mathcal{S}}^{r,auth}(\lambda)$, is quite similar to that above. The main difference compared to the previous experiment is that the adversary \mathcal{A} tries to make some legitimate tag to authenticate the reader. As π and \mathcal{T}_{ID} have no matching conversation, \mathcal{A} computes at least one message that makes the tag to authenticate the reader.

An RFID scheme \mathcal{S} achieves *reader authentication* if the advantage of \mathcal{A} , $Adv_{\mathcal{A},\mathcal{S}}^{r,auth}$, is negligible, for any strong adversary \mathcal{A} ($Adv_{\mathcal{A},\mathcal{S}}^{r,auth}$ is defined as above, by using $RFID_{\mathcal{A},\mathcal{S}}^{r,auth}(\lambda)$ instead of $RFID_{\mathcal{A},\mathcal{S}}^{t,auth}(\lambda)$).

Privacy. *Privacy* for RFID systems [2] captures anonymity and untraceability. It basically means that an adversary cannot learn anything new from intercepting the communication between a tag and the reader. To model this, the concept of a *blinder* was introduced in [2].

A *blinder* for an adversary \mathcal{A} that belongs to some class V of adversaries is a PPT algorithm \mathcal{B} that simulates the *Launch*, *SendReader*, *SendTag*, and *Result* oracles for \mathcal{A} , without having access to the corresponding secrets. Moreover, it looks passively at the communication between \mathcal{A} and the other oracles allowed to it by the class V (that is, \mathcal{B} gets exactly the same information as \mathcal{A} when querying these oracles).

When the adversary \mathcal{A} interacts with the RFID scheme by means of a blinder \mathcal{B} , we say that \mathcal{A} is *blinded by \mathcal{B}* and denote this by $\mathcal{A}^{\mathcal{B}}$.

Given an adversary \mathcal{A} , define the experiment (privacy game):

Experiment $RFID_{\mathcal{A},\mathcal{S}}^{prv-0}(\lambda)$

- 1: Set up the reader;
- 2: \mathcal{A} gets the public key pk ;
- 3: \mathcal{A} queries the oracles;
- 4: \mathcal{A} gets the secret table of the *DrawTag* oracle;
- 5: \mathcal{A} outputs a bit b' ;
- 6: Return b' .

In the same way, by replacing “ \mathcal{A} ” with “ $\mathcal{A}^{\mathcal{B}}$ ”, we define the experiment $RFID_{\mathcal{A},\mathcal{S},\mathcal{B}}^{prv-1}(\lambda)$. Now, the *advantage* of \mathcal{A} blinded by \mathcal{B} is

$$Adv_{\mathcal{A},\mathcal{S},\mathcal{B}}^{prv}(\lambda) = |P(RFID_{\mathcal{A},\mathcal{S}}^{prv-0}(\lambda) = 1) - P(RFID_{\mathcal{A},\mathcal{S},\mathcal{B}}^{prv-1}(\lambda) = 1)|.$$

An RFID scheme is private for a class V of adversaries if, for any $\mathcal{A} \in V$, there exists a blinder \mathcal{B} such that $Adv_{\mathcal{A}, \mathcal{S}, \mathcal{B}}^{prv}(\lambda)$ is negligible.

We thus obtain eight concepts of privacy: *strong privacy*, *narrow strong privacy*, *destructive privacy*, and so on.

4 Destructive privacy and reader-first authentication

An interesting question that arises when designing mutual authentication RFID schemes is whether the tag or the reader should be authenticated first. We have thus two approaches: *tag-first* and *reader-first authentication*, respectively [3]. The *tag-first authentication* has some advantage with respect to desynchronization: the tag computes its new state and sends information about it to the reader. However, the tag state is updated only when the reader authenticates the tag and confirms the new state to the tag. The disadvantage of this approach is that the tag should provide some information to the reader before it is confident of the reader’s identity.

The *reader-first authentication* might enhance the tag privacy because the tag gives private information to the reader when it is confident of its identity. This also might help preventing adversaries from tracking tags. Another advantage is when the tag is designed only for a limited number of authentications. In such a case, the reader-first approach prevents a form of the denial of service attack that would “consume” all the tag’s authentication answers.

In this section, we address the problem to construct a destructive private and mutual authentication RFID scheme in Vaudenay’s model with temporary state disclosure. For mutual authentication, we follow the reader-first approach and endow all tags with PUFs.

To describe our scheme, let us assume that λ is a security parameter, $\ell_1(\lambda)$ and $\ell_2(\lambda)$ are two polynomials, and $F = (F_K)_{K \in \mathcal{K}}$ is a pseudo-random function, where $F_K : \{0, 1\}^{2\ell_1(\lambda)+1} \rightarrow \{0, 1\}^{\ell_1(\lambda)}$ for all $K \in \mathcal{K}_\lambda$. Each tag is equipped with a (unique) PUF $P : \{0, 1\}^{p(\lambda)} \rightarrow \mathcal{K}_\lambda$ and has the capacity to compute F , where $p(\lambda)$ is a polynomial. The internal state of the tag consists of a pair (s, x) , where $s \in \{0, 1\}^{p(\lambda)}$ is randomly chosen as a seed to evaluate P , and $x \in \{0, 1\}^{\ell_1(\lambda)}$ is a

random string that gets incremented after each protocol instance. The reader maintains a database DB with entries for all legitimate tags. Each entry is a vector (ID, K) , where ID is the tag's identity and $K = P(s)$, where P is the tag's PUF, and (s, x) is its state.

	Reader (DB, F)	Tag (P, s, F, x)
1		$x = x + 1, K = P(s)$ $u = F_K(0, 0, x), z = F_K(0, u, 0)$ erase K, u, z
2	If $\exists (ID, K) \in DB$ s.t. $z = F_K(0, u, 0)$ then $v \leftarrow \{0, 1\}^{\ell_1(\lambda)}$, $w = F_K(1, v, u)$ else $v \leftarrow \{0, 1\}^{\ell_1(\lambda)}$ $w \leftarrow \{0, 1\}^{\ell_1(\lambda)}$	$\xrightarrow{v, w}$
3		$K = P(s), u = F_K(0, 0, x)$ If $w = F_K(1, v, u)$ then $w' = F_K(1, v, u + 1)$ else $w' = F_K(1, v + 1, u)$ erase K, v, w, w'
4	If $w' = F_K(1, v, u + 1)$ then output ID else output \perp	$\xleftarrow{w'}$

Figure 1. Destructive private and reader-first authentication

The mutual authentication protocol is given in Figure 1. As we can see, the tag initially increments x and computes $K = P(s)$, $u = F_K(0, 0, x)$, $z = F_K(0, u, 0)$. The tuple (u, z) is then sent to the reader. The reader checks its database for a tuple (ID, K) such that $z = F_K(0, u, 0)$. When the reader finds out the right value, it prepares the answer for the tag by generating a random $v \leftarrow \{0, 1\}^{\ell_1(\lambda)}$ and computes $w = F_K(1, v, u)$. If no such entry is found, then the reader chooses both v and w as random values. The tag evaluates the PUF,

checks the value w received from the reader, takes a decision, and prepares the answer for the reader. On receiving the tag's answer, the reader checks it and takes a decision.

Theorem 4.1. *The RFID scheme in Figure 1 is correct.*

Proof. Assuming that a tag \mathcal{T}_{ID} is legitimate, the reader's database contains an entry (ID, K) , where $K = P(s)$, (s, x) is the tag's state, and P is its PUF.

When the reader receives (u, z) from the tag \mathcal{T}_{ID} , exactly the equality $z = F_K(0, u, 0)$ holds with overwhelming probability (we use the same notation as in Figure 1).

If the reader has found the tag in its database (i.e., identified it), the equality $w = F_K(1, v, u)'$ on tag's side holds with overwhelming probability. This means that the tag authenticates the reader. In such a case, the equality $w' = F_K(1, v, u + 1)$ holds with overwhelming probability, meaning that the reader authenticates the tag.

As a final remark, if the tag does not authenticate the reader, then the reader will authenticate the tag with negligible probability. \square

We will focus now on the security of our RFID scheme.

Theorem 4.2. *The RFID scheme in Figure 1 achieves tag authentication in Vaudenay's model with temporary state disclosure, provided that F is a PRF and the tags are endowed with ideal PUFs.*

Proof. Assume that the scheme does not achieve tag authentication, and let \mathcal{A} be an adversary that has non-negligible advantage over the scheme, with respect to the tag authentication property. We will show that there exists a PPT algorithm \mathcal{A}' that can break the pseudo-randomness property of the function F .

The main idea is the next one. Let \mathcal{C} be a challenger for the pseudo-randomness security game of the function F . The adversary \mathcal{A}' will play the role of challenger for \mathcal{A} . Thus, \mathcal{A}' guesses the identity ID^* of an uncorrupted legitimate tag that has no matching conversation with the reader, but \mathcal{A} can make the reader authenticate it with a non-negligible probability (recall that there is a polynomial number $t(\lambda)$ of tags). Then, it creates the tag \mathcal{T}_{ID^*} with the help of \mathcal{C} . Specifically, \mathcal{A}'

will not associate any key with the identity ID^* . However, the function chosen by \mathcal{C} will be the one used to do all the calculations of this tag. This function is either $F_{K^*} \leftarrow F$ for some K^* , or a random function. \mathcal{T}_{ID^*} will be regarded by \mathcal{A} as a legitimate tag. The adversary \mathcal{A}' does not know this function but, with the help of \mathcal{A} , he will try to distinguish between the two cases with non-negligible probability.

The details on \mathcal{A}' are as follows (λ is a security parameter):

1. The challenger \mathcal{C} chooses uniformly at random $F_{K^*} \leftarrow F$ for some K^* , or a random function. Let us denote it by f ;
2. \mathcal{A}' plays the role of challenger for \mathcal{A} . It will run the reader and all tags created by \mathcal{A} , answering all \mathcal{A} 's oracle queries. Therefore, using $SetupR(\lambda)$, it generates a triple (pk, sk, DB) , gives the public key pk to \mathcal{A} , and keeps the private key sk .

\mathcal{A}' will maintain a list of tag entries $\mathcal{A}'_{ListTags}$ similar to $ListTags$ (see Section 3) but with the difference that each entry in this list also includes the current state of the tag as well as a special field designated to store the “key generated by the tag’s internal PUF”. The legitimate entries in this list define the reader’s database DB . Initially, $\mathcal{A}'_{ListTags}$ is empty;

3. \mathcal{A}' guesses the tag identity ID^* that \mathcal{A} will authenticate to reader (please see the discussion above);
4. \mathcal{A}' will simulate for \mathcal{A} all the corresponding oracles in a straightforward manner, but with the following modifications:

- (a) $CreateTag^b(ID)$: If \mathcal{T}_{ID} was already created, then \mathcal{A}' does nothing. If \mathcal{T}_{ID} was not created and $ID \neq ID^*$, then \mathcal{A}' randomly chooses $K \in \{0, 1\}^\lambda$ and $x \in \{0, 1\}^{\ell_1(\lambda)}$ and records a corresponding entry into $\mathcal{A}'_{ListTags}$ (K plays the role of the key generated by the tag’s internal PUF). Thus, \mathcal{T}_{ID} has just been created. If \mathcal{T}_{ID} was not created and $ID = ID^*$, then \mathcal{A}' records $(ID^*, ?, x)$ into $\mathcal{A}'_{ListTags}$, where $x \leftarrow \{0, 1\}^{\ell_1(\lambda)}$. The meaning of “?” is that this field should have contained a key for F . However, \mathcal{A}' does not even know if \mathcal{C} chose a function from F (so a key) or a random function. However, \mathcal{A}' does not need to know this because it can answer all \mathcal{A} 's queries regarding ID^* with the help of \mathcal{C} .

As the tags are endowed with ideal PUFs and the keys are uniformly at random chosen by \mathcal{A}' , including the function chosen by \mathcal{C} , \mathcal{A}' implements correctly the functionality of all tags (including \mathcal{T}_{ID^*});

- (b) *DrawTag* and *Free*: \mathcal{A}' knows the list of all tags created by \mathcal{A} , and updates it correspondingly whenever \mathcal{A} draws or frees some tag;
- (c) *Launch()*: \mathcal{A}' launches a new protocol instance whenever \mathcal{A} asks for it;
- (d) *SendTag*($\emptyset, vtag$): This is the first message $vtag$ sends in a protocol instance. If the tag referred by $vtag$ is ID^* , then \mathcal{A}' will increment x and then query \mathcal{C} for $(0, 0, x)$, which will become u , and for $(0, u, 0)$, which will become z . If u, z are \mathcal{C} 's responses, then \mathcal{A}' answers with (u, z) .

If $vtag$ refers to some $ID \neq ID^*$, then \mathcal{A}' can prepare the answer because it knows the corresponding key for ID ;

- (e) *SendReader*($((u, z), \pi)$): Assume the reader (run by \mathcal{A}') has received (u, z) in the protocol instance π from a tag identified by $vtag$ (in other words, $(u, z) \leftarrow \text{SendTag}(\emptyset, vtag)$).

If $vtag$ refers to some tag ID such that $(ID, K) \in DB$ for some K , then the reader (run by \mathcal{A}') can compute the answer according to the protocol.

If $vtag$ refers to ID^* , then the reader (run by \mathcal{A}') can compute the answer according to the protocol by querying \mathcal{C} (recall that \mathcal{T}_{ID^*} is regarded by \mathcal{A} as a legitimate tag).

If $vtag$ refers to some ID for which no entry can be found in DB , then the answer (v, w) is randomly chosen;

- (f) *SendTag*($((v, w), vtag)$): If the tag referred by $vtag$ is ID^* , then \mathcal{A}' queries \mathcal{C} for $(1, v, u)$ and then compares the answer with w . If they match, the tag outputs OK ; otherwise, it outputs \perp . In the first case, \mathcal{A}' queries \mathcal{C} for $(1, v, u + 1)$ to get w' ; in the second case, it queries \mathcal{C} for $(1, v + 1, u)$. If $vtag$ refers to some $ID \neq ID^*$ that has associated a pair (K, x) , then \mathcal{A}' can compute by itself w' (according to the protocol). In all cases, the oracle returns w' ;

- (g) $Result(\pi)$: \mathcal{A}' can infer the decision of the reader in the last step of π because it can obtain the value $F_K(1, v, u)$ for all tags (either it can compute it or query \mathcal{C} for it). Therefore, \mathcal{A}' can simulate $Result(\pi)$ according to its definition;
- (h) $Corrupt(vtag)$: If the tag referred by $vtag$ is different from ID^* , then \mathcal{A}' returns its current state; otherwise, it aborts.

If \mathcal{A}' sees that \mathcal{A} could make the reader authenticate \mathcal{T}_{ID^*} without corrupting it and without any matching conversation between tag and reader, it answers to \mathcal{C} that f is PRF; otherwise, f is random. In the first case,

$$P(1 \leftarrow (\mathcal{A}')^{F_K}(1^\lambda) : K \leftarrow \mathcal{K}_\lambda) = P(RFID_{\mathcal{A}, \mathcal{S}}^{t, auth}(\lambda) = 1).$$

In the second case,

$$P(1 \leftarrow \mathcal{A}^f(1^\lambda) : f \leftarrow U_\lambda) = \eta(\lambda)$$

for some negligible function $\eta(\lambda)$. That is because \mathcal{A} does not play the real tag authentication game, the function implemented by \mathcal{T}_{ID^*} is random, and the tag does not have any matching conversation with the reader. So, the reader (simulated by \mathcal{A}') should authenticate the tag on behalf of a random message (u, z) sent by \mathcal{A} to reader, message that is verified for correctness by \mathcal{C} .

Therefore,

$$Adv_{\mathcal{A}, F}^{prf}(\lambda) = |P(RFID_{\mathcal{A}, \mathcal{S}}^{t, auth}(\lambda) = 1) - \eta(\lambda)|.$$

If we assume now that \mathcal{A} has a non-negligible probability to make the reader authenticate the tag \mathcal{T}_{ID^*} , then \mathcal{A}' will have a non-negligible advantage against F ; this contradicts the fact that F is a pseudo-random function. \square

As with respect to the reader authentication property, we have the following result.

Theorem 4.3. *The RFID scheme in Figure 1 achieves reader authentication in Vaudenay's model with temporary state disclosure, provided that F is a PRF and the tags are endowed with ideal PUFs.*

Proof. Assume that our scheme does not achieve reader authentication, and let \mathcal{A} be an adversary that has a non-negligible advantage over the scheme, with respect to the reader authentication property. We will show that there exists a PPT algorithm \mathcal{A}' that can break the pseudo-randomness property of the function F .

The main idea is somewhat similar to the one in the Theorem 4.2. Let \mathcal{C} be a challenger for the pseudo-randomness property of the function F . The adversary \mathcal{A}' will play the role of a challenger for \mathcal{A} . First, \mathcal{A}' guesses the identity ID^* of an uncorrupted legitimate tag that has no matching conversation with the reader, but \mathcal{A} can make the tag authenticate the reader with a non-negligible probability (recall that there is a polynomial number $t(\lambda)$ of tags). Then, it creates the tag \mathcal{T}_{ID^*} with the help of \mathcal{C} , exactly as in the proof of Theorem 4.2. This tag will be regarded by \mathcal{A} as a legitimate one.

The description of \mathcal{A}' is very similar to the one in the proof of Theorem 4.2, so we will focus on the differences between them (λ denotes a security parameter):

1. The challenger \mathcal{C} chooses uniformly at random $F_{K^*} \leftarrow F$ for some K^* , or a random function. Let us denote it by f ;
2. \mathcal{A}' plays the role of a challenger for \mathcal{A} . It will run the reader and all tags created by \mathcal{A} , answering all \mathcal{A} 's oracle queries. Therefore, using $SetupR(\lambda)$ it generates a triple (pk, sk, DB) , gives the public key pk to \mathcal{A} , and keeps the private key sk .

\mathcal{A}' will maintain a list of tag entries $\mathcal{A}'_{ListTags}$ exactly as in the proof of Theorem 4.2;

3. \mathcal{A}' guesses the tag identity ID^* that authenticates \mathcal{A} as a valid reader;
4. \mathcal{A}' will simulate for \mathcal{A} all the corresponding oracles exactly as in the proof of Theorem 4.2.

The advantage of \mathcal{A}' against the PRF F is computed as in the proof of Theorem 4.2. Therefore, the assumption that \mathcal{A} has a non-negligible probability to make \mathcal{T}_{ID^*} authenticate it as a valid reader contradicts the pseudo-randomness of the function F . \square

By using the *sequence-of-games* approach [19], we will prove that our protocol reaches destructive privacy. With this approach, a sequence of games (probabilistic experiments) is defined. The initial game is the original privacy game with respect to a given adversary. The transition from one game G_i to another one G_{i+1} is done by indistinguishability in our case. This means that a probability distribution in G_i is replaced by another one that is indistinguishable from the previous one. In this way, the difference between the probabilities the adversary wins G_i and G_{i+1} , is negligible.

Theorem 4.4. *The RFID scheme in Figure 1 achieves destructive privacy in Vaudenay’s model with temporary state disclosure, provided that F is a PRF and the tags are endowed with ideal PUFs.*

Proof. Let \mathcal{A} be a destructive adversary against our RFID scheme denoted \mathcal{S} . We will show that there is a blinder \mathcal{B} such that $Adv_{\mathcal{A},\mathcal{S},\mathcal{B}}^{prv}(\lambda)$ is negligible. The blinder \mathcal{B} that we construct, which has to answer to the oracles *Launch*, *SendReader*, *SendTag*, and *Result* without knowing any secret information, works as follows:

- *Launch*(\emptyset): returns a unique identifier π for a new protocol instance;
- *SendTag*($\emptyset, vtag$): returns (u, z) , where $u, z \leftarrow \{0, 1\}^{\ell_1(\lambda)}$;
- *SendReader*($((u, z), \pi)$): returns (v, w) , where $v, w \leftarrow \{0, 1\}^{\ell_1(\lambda)}$;
- *SendTag*($((v, w), vtag)$): returns $w' \leftarrow \{0, 1\}^{\ell_1(\lambda)}$;
- *SendReader*((w', π)): the blinder does not do anything because, in this case, the reader does not answer;
- *Result*(π): if the session π does not exist or exists but is not completed, the blinder outputs \perp . If π has been issued by the *Launch*(\emptyset) oracle and a protocol transcript $tr_\pi = ((u, z), (v, w), w')$ has been generated by

- $(u, z) \leftarrow \text{SendTag}(\emptyset, vtag)$,
- $(v, w) \leftarrow \text{SendReader}((u, z), \pi)$,
- $w' \leftarrow \text{SendTag}((v, w), vtag)$, and
- $\text{SendReader}(w', \pi)$,

where $vtag$ refers to some legitimate tag, the blinder outputs 1; otherwise, outputs 0 (remark that the blinder sees what \mathcal{A} sees and, therefore, it knows whether $vtag$ refers to some legitimate tag or not).

We further prove that $Adv_{\mathcal{A},\mathcal{S},\mathcal{B}}^{prv}(\lambda)$ is negligible. To this we define a sequence of games G_0, \dots, G_7 , where G_0 is the experiment $RFID_{\mathcal{A},\mathcal{S}}^{prv-0}$ and G_{i+1} is obtained from G_i as described below, for all $0 \leq i < 7$. By $P(G_i)$ we denote the probability the adversary \mathcal{A} wins the game G_i .

Game G_1 : This is identical to G_0 except that the game challenger will not use the PRF keys generated by PUFs to answer the adversary's oracle queries, but randomly generated keys, one for each tag created by adversary. Of course, the game challenger must maintain a secret table with the association between each tag and this new secret key. From the adversary's point of view, this means that the probability distribution given by each tag's PUF (in G_0) is replaced by the uniform probability distribution (in G_1). As the PUFs are ideal, the two distributions are indistinguishable. Taking into account that there are a polynomial number of tags, it must be the case that $|P(G_0) - P(G_1)|$ is negligible.

Game G_2 : We replace in G_1 the oracle $Result$ by $Result_{\mathcal{B}}$, which is the simulation of $Result$ by the blinder \mathcal{B} (please, see the definition \mathcal{B}). Denote by G_2 the game such obtained. We prove that $P(G_1) = P(G_2)$.

Recall first that in game G_1 the tags are still endowed with PUFs, but their secret PRF keys are not computed by PUFs. They are randomly generated by the game challenger that maintains a secret table with the key associated to each tag. In this way, the *Corrupt* oracle will never reveal the secret key, but it destroys the tag when queried.

If \mathcal{A} queries $Result$ or $Result_{\mathcal{B}}$ for a protocol session that does not exist or is incomplete, both oracles return \perp . Therefore, let us assume that these oracles are queried on a complete protocol session π . In this case, we will show that $Result(\pi) = 1$ if and only if $Result_{\mathcal{B}}(\pi) = 1$.

Assume $Result(\pi) = 1$. Then, there is a transcript $tr_{\pi} = ((u, z), (v, w), w')$ defined by a sequence of oracle queries

- $(u, z) \leftarrow SendTag(\emptyset, vtag)$
- $(v, w) \leftarrow SendReader((u, z), \pi)$

- $w' \leftarrow \text{SendTag}((v, w), vtag)$
- and $\text{SendReader}(w', \pi)$

such that $vtag$ refers to some tag \mathcal{T}_{ID} whose state is (s, x) and secret key is K , $u = F_K(0, 0, x)$, $z = F_K(0, 0, u)$, and (ID, K) is in the reader's database (that is, \mathcal{T}_{ID} is legitimate). All these facts show that $\text{Result}_{\mathcal{B}}(\pi) = 1$ (recall that the blinder \mathcal{B} sees what \mathcal{A} sees and, therefore, it knows whether $vtag$ refers to some legitimate tag or not).

The inverse implication is a bit more elaborate. Assume that $\text{Result}_{\mathcal{B}}(\pi) = 1$. This means that there is a transcript $tr_{\pi} = ((u, z), (v, w), w')$ defined by a sequence of oracle queries as those above and the tag \mathcal{T}_{ID} referred by $vtag$ is legitimate. Assume that the tag's key is K and its state is (s, x) , and in DB there is a record (ID, K) . Because the oracles SendReader and SendTag are the real ones (and not simulated by blinder), the reader finds a record such that $z = F_K(0, 0, u)$. Therefore, w must be of the form $F_K(0, v, u)$, and this value will match $F_K(0, v, u)$ computed by tag. Therefore, the tag authenticates the reader and replies by $w' = F_K(1, v, u)$. But then, the reader will successfully check the equality between w and $F_K(1, v, u)$ (computed by itself) and, therefore, authenticates the tag. As a conclusion, $\text{Result}(\pi) = 1$.

This shows that $P(G_1) = P(G_2)$.

Game G_3 : This game is identical to G_2 , except that the $\text{Launch}()$ oracle is simulated according to the blinder description. No difference is encountered between the two games and, therefore, $P(G_2) = P(G_3)$.

Game G_4 : This is identical to G_3 except that the $\text{SendTag}(\emptyset, vtag)$ oracle is simulated according to the blinder description. By doing this, the probability distribution $\{(u, z) \mid u = F_K(0, 0, x), z = F_K(0, u, 0)\}$ is replaced by $\{(u, z) \mid u, z \leftarrow \{0, 1\}^{\ell_1(\lambda)}\}$.

As F is a PRF, $|P(G_3) - P(G_4)|$ is negligible. The proof is quite straightforward. The main idea is as follows. Assume that an adversary \mathcal{A} can distinguish with a non-negligible probability between G_3 and G_4 . Define an adversary \mathcal{A}' for PRF that uses \mathcal{A} as a subroutine and send $(0, u, 0)$ as a challenge. When the PRF challenger returns, with equal probability, either $z = F_K(0, u, 0)$ or $z \leftarrow \{0, 1\}^{\ell_1(\lambda)}$, \mathcal{A}' sends this

value to \mathcal{A} . The probability \mathcal{A}' guesses between the two possibilities for z is exactly the probability \mathcal{A} distinguishes between the two games.

Game G_5 : This game is identical to G_4 , except that the oracle $SendReader((u, z), \pi)$ is simulated according to the blinder description. That is, for each tag \mathcal{T}_{ID} whose secret key is K and current state is (s, x) , one of the two probability distributions

$$\begin{cases} \{(u, z, v, w) \mid u, v, z \leftarrow \{0, 1\}^{\ell_1(\lambda)}, w = F_K(1, v, u)\}, \\ \{(u, z, v, w) \mid u, v, z, w \leftarrow \{0, 1\}^{\ell_1(\lambda)}\}, \end{cases}$$

is replaced by $\{(u, z, v, w) \mid u, v, z, w \leftarrow \{0, 1\}^{\ell_1(\lambda)}\}$.

As F is a PRF, and the key K was chosen at random, it must be the case that $|P(G_4) - P(G_5)|$ is negligible. The proof is by contradiction, and it is quite similar to the proof that establishes the transition from G_3 to G_4 .

Game G_6 : This game is identical to G_5 , except that the oracle $SendTag((v, w), vtag)$ is simulated by blinder. That is, for each tag \mathcal{T}_{ID} , one of the two probability distributions

$$\begin{cases} \{(u, z, v, w, w') \mid u, v, z, w \leftarrow \{0, 1\}^{\ell_1(\lambda)}, w' = F_K(1, v, u + 1)\}, \\ \{(u, z, v, w, w') \mid u, v, z, w \leftarrow \{0, 1\}^{\ell_1(\lambda)}, w' = F_K(1, v + 1, u)\}, \end{cases}$$

is replaced by $\{(u, z, v, w, w') \mid u, v, z, w, w' \leftarrow \{0, 1\}^{\ell_1(\lambda)}\}$.

As F is a PRF, and the key K was chosen at random, it must be the case that $|P(G_5) - P(G_6)|$ is negligible. The proof is by contradiction, and it is quite similar to the proof in Game G_5 . Therefore, it is omitted.

Game G_7 : This is identical to G_6 , except that $SendReader(w', \pi)$ is simulated by blinder. However, this does not change the probability distribution from G_6 . Therefore, $P(G_6) = P(G_7)$.

Now, we show that G_7 is in fact $RFID_{\mathcal{A}, S, \mathcal{B}}^{prv-1}$. The blinded adversary $\mathcal{A}^{\mathcal{B}}$ sees each tag as a standard PUF tag, although random secret keys are used instead of the keys generated by PUFs. The oracles $CreateTag$, $Draw$, $Free$, and $Corrupt$ that can be queried directly by \mathcal{A} do not use the keys generated by PUFs in order to answer the adversary's queries (in fact, they do not use any secret key). The answer to the other oracles is simulated by a blinder that does not use the secret keys either. Therefore, G_7 is indeed $RFID_{\mathcal{A}, S, \mathcal{B}}^{prv-1}$.

Now, remark that $P_{\mathcal{A}}(G_0) = P(\text{RFID}_{\mathcal{A},\mathcal{S}}^{\text{prv}-0}(\lambda) = 1)$ and $P_{\mathcal{A}}(G_7) = P(\text{RFID}_{\mathcal{A},\mathcal{S},\mathcal{B}}^{\text{prv}-1}(\lambda) = 1)$. Combining all the probabilities $P(G_i)$ together, we obtain that $\text{Adv}_{\mathcal{A},\mathcal{S},\mathcal{B}}^{\text{prv}}(\lambda)$ is negligible and, therefore, our protocol achieves destructive privacy. \square

5 Narrow destructive privacy and reader-first authentication

With little effort, we can design a similar scheme that achieves narrow destructive privacy and reader-first authentication RFID in Vaudenay’s model with temporary state disclosure. The mutual authentication protocol of this new RFID scheme is presented in Figure 2; all the other elements are as in Section 4, except that F_K is a function from $\{0, 1\}^{\ell_1(\lambda)+2}$ to $\{0, 1\}^{\ell_2(\lambda)}$ and t is polynomial in the security parameter.

As one can see, there is no random generator on tag. Because of this, the synchronization between tag and reader can be lost. The only thing we can do is to check (on the reader side) for a polynomial bounded desynchronization. Due to this, the scheme can be at most narrow destructive private: if an adversary desynchronizes the tag and reader sufficiently enough (for more than t steps), then it will be able to distinguish the real privacy game from the blinded one by means of the *Result* oracle. Roughly speaking, this is because, in the real privacy game, the *Result* oracle returns 0 (when the tag and reader are desynchronized for more than t steps), while in the blinded privacy game, it returns 1. We, therefore, have the following result.

Theorem 5.1. *The RFID scheme in Figure 2 achieves mutual authentication and narrow destructive privacy in Vaudenay’s model with temporary state disclosure, provided that F is a PRF and the tags are endowed with ideal PUFs.*

Proof. It is straightforward to see that the proof follows a similar line to the proofs of Theorems 4.2 and 4.3 for mutual authentication, and Theorem 4.4 for narrow destructive privacy. Remark that for privacy, the *Result* oracle is not used. \square

	Reader (DB, F)	Tag (P, s, F, x)
1		\xleftarrow{z} $K = P(s), z = F_K(0, 0, x)$ erase K, z $x = x + 1$
2	If $\exists(ID, K, x) \in DB$ and $0 \leq i < t$ s.t. $z = F_K(0, 0, x + i)$ then $x = x + i$ $w = F_K(0, 1, x + 1)$ else $w \leftarrow \{0, 1\}^{\ell_2(\lambda)}$	\xrightarrow{w}
3		$K = P(s)$ If $w \neq F_K(0, 1, x)$ then $w' = F_K(1, 1, x)$ else $w' = F_K(1, 0, x)$ $\xleftarrow{w'}$ erase K, w, w'
	If $w' = F_K(1, 1, x + 1)$ then output $ID, x = x + 1$ else output \perp	

Figure 2. Narrow destructive private and reader-first authentication

It is good to remark that our RFID scheme in Figure 2 also provides an appropriate practical solution to the narrow destructive privacy in the plain Vaudenay’s model, where the existing solution is based on random oracles [1], [2] .

A few more words on desynchronization are in order. If we look to the protocol in Figure 2, we remark that the desynchronization is a result of the fact that the tag and reader share a common variable x that is updated by tag before authenticating the reader. This allows an adversary to query a tag for more than t times and, therefore, to desynchronize the tag and the reader.

To prevent desynchronization between reader and tag in reader-first authentication RFID schemes, the tag should update the shared permanent variables after authenticating the reader, and not before.

6 Conclusions

Modern applications of RFID systems ask for advanced security and privacy properties. For instance, tag destruction under corruption is an important requirement when the tag is used for access control. Likewise, the disclosure of temporary state under tag corruption is a serious threat in practice. Reader-first authentication [3] assures that the tag will give its private data only when it authenticates the reader. Therefore, tag tracking and data theft is prevented when the reader is fake. All these together mean that we need RFID schemes that provide destructive privacy and reader-first authentication under corruption with temporary state disclosure.

The aim of this paper is to propose two RFID schemes that fill this gap. The first one is destructive private and the second one is narrow destructive private. Both of them assure reader-first authentication, are practical, and efficient. Also, both schemes avoid random number generators on tags. As (narrow) destructive privacy cannot be achieved with ordinary tags, we have used PUFs as secure hardware containers for the secret key of tags. Detailed security and privacy proofs are provided for our schemes.

References

- [1] S. Vaudenay, “On privacy models for RFID,” in *Proceedings of the Advances in Cryptology 13th International Conference on Theory and Application of Cryptology and Information Security*, ser. ASIACRYPT’07. Berlin, Heidelberg: Springer-Verlag, 2007, pp. 68–87.
- [2] R.-I. Païse and S. Vaudenay, “Mutual authentication in RFID: Security and privacy,” in *Proceedings of the 2008 ACM Symposium on Information, Computer and Communications Security*, ser. ASIACCS ’08. New York, USA: ACM, 2008, pp. 292–299.
- [3] R. Peeters, J. Hermans, and J. Fan, “IBIHOP: Proper privacy preserving mutual RFID authentication,” in *The 2013 Workshop on Radio Frequency Identification/Internet of Things Security (RFIDsec’13 Asia)*, ser. Cryptology and Information Security Series, C. Ma and J. Weng, Eds., vol. 11. IOS Press, 2013.

- [4] C. Hristea and F. L. Țiplea, “Destructive privacy and mutual authentication in Vaudenay’s RFID model,” Cryptology ePrint Archive, Report 2019/073, 2019, <https://eprint.iacr.org/2019/073>.
- [5] S. Kardaş, S. Çelik, M. Yildiz, and A. Levi, “PUF-enhanced of-line RFID security and privacy,” *J. Netw. Comput. Appl.*, vol. 35, no. 6, pp. 2059–2067, Nov. 2012.
- [6] M. Akgün and M. U. Çaglayan, “Providing destructive privacy and scalability in RFID systems using PUFs,” *Ad Hoc Netw.*, vol. 32, no. C, pp. 32–42, Sep. 2015.
- [7] F. Armknecht, A.-R. Sadeghi, A. Scafuro, I. Visconti, and C. Wachsmann, “Impossibility results for RFID privacy notions,” in *Transactions on Computational Science XI*, M. L. Gavrilova, C. J. K. Tan, and E. D. Moreno, Eds. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 39–63.
- [8] A.-R. Sadeghi, I. Visconti, and C. Wachsmann, “PUF-enhanced RFID security and privacy,” in *Workshop on secure component and system identification (SECSI)*, vol. 110, 2010.
- [9] —, *Enhancing RFID Security and Privacy by Physically Unclonable Functions*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 281–305.
- [10] M. Sipser, *Introduction to the Theory of Computation*. Cengage Learning, 2012.
- [11] K. Finkenzeller, *RFID Handbook: Fundamentals and Applications in Contactless Smart Cards and Identification*, 3rd ed. Wiley Publishing, 2010.
- [12] Y. Li, H. R. Deng, and E. Bertino, *RFID Security and Privacy*, ser. Synthesis Lectures on Information Security, Privacy, and Trust. Morgan & Claypool Publishers, 2013.
- [13] A. Juels and S. A. Weis, “Defining strong privacy for RFID,” *ACM Trans. Inf. Syst. Secur.*, vol. 13, no. 1, pp. 7:1–7:23, Nov. 2009.
- [14] S. Canard, I. Coisel, J. Etrog, and M. Girault, “Privacy-preserving RFID systems: Model and constructions,” <https://eprint.iacr.org/2010/405.pdf>, 2010.
- [15] R. H. Deng, Y. Li, M. Yung, and Y. Zhao, “A new framework for RFID privacy,” in *Proceedings of the 15th European Conference*

- on Research in Computer Security*, ser. ESORICS'10. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 1–18.
- [16] J.-M. Bohli and A. Pashalidis, “Relations among privacy notions,” *ACM Trans. Inf. Syst. Secur.*, vol. 14, no. 1, pp. 4:1–4:24, Jun. 2011.
- [17] J. Hermans, F. Pashalidis, Andreasand Vercauteren, and B. Preneel, “A new RFID privacy model,” in *Computer Security – ESORICS 2011*, V. Atluri and C. Diaz, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 568–587.
- [18] J. Hermans, R. Peeters, and B. Preneel, “Proper RFID privacy: Model and protocols,” *IEEE Transactions on Mobile Computing*, vol. 13, no. 12, pp. 2888–2902, Dec 2014.
- [19] V. Shoup, “Sequences of games: A tool for taming complexity in security proofs,” 2004.

Ferucio Laurențiu Țiplea, Cristian Hristea,
Rodica Bulai

Received April 11, 2022
Revised April 28, 2022
Accepted April 28, 2022

Ferucio Laurențiu Țiplea
Department of Computer Science
“Alexandru Ioan Cuza” University of Iași
Iași, Romania
E-mail: ferucio.tiplea@uaic.ro

Cristian Hristea
Simion Stoilow Institute of Mathematics of the Romanian Academy
Bucharest, Romania
E-mail: cristi.hristea@gmail.com

Rodica Bulai
Faculty of Computers, Informatics and Microelectronics
Technical University of Moldova
Chisinau, Republic of Moldova
E-mail: rodica.bulai@ati.utm.md

Kurtosis-Based Feature Selection Method using Symmetric Uncertainty to Predict the Air Quality Index

Usharani Bhimavarapu, M. Sreedevi

Abstract

Feature selection is vital in data pre-processing in machine learning, and it is prominent in datasets with many features. Feature selection analyses the relevant, irrelevant, and redundant features in the dataset. Feature selection removes the irrelevant features, which improves both the accuracy and prediction performance. The significant advantages of reducing the number of features from the dataset are reducing the training time, reducing overfitting, decreasing the curse of dimensionality, and simplifying the prediction model. The filter feature selection techniques can handle the issues with the high number of features, and this paper uses the symmetric uncertainty coefficient to verify the relevance of the independent features. In this paper, a new feature selection method named as kurtosis-based feature selection has been proposed to select the relevant features which affect the air pollution. Kurtosis-based feature selection is compared with seven filter feature selection techniques on air pollution dataset and validated the performance of the proposed algorithm. It has been observed that the kurtosis-based feature selection extracts only PM2.5 as the key feature and has been compared to the accuracy of the five existing methods. The experimental results illustrate that the kurtosis-based feature selection algorithm reduces the original feature set up to 91.66%, but the existing filter feature selection techniques reduce the feature set to only 50%.

Keywords: Air Pollution, Air quality index, Correlation coefficient, Feature selection, Filter techniques, Symmetric uncertainty.

1 Introduction

The features in the dataset may be repeated and noisy, and these repeated reduce the learning model's performance. This paper explains how to choose the relevant features to predict the air quality index and ignore the irrelevant features. The main advantages of feature selection (FS) are improving the prediction performance by removing the irrelevant and redundant features and reducing the computational cost [1].

There are two types of FS techniques: classifier independent (filter) and classifier dependent (wrapper, embedded) [2]. Filter FS techniques give the grade for each feature and select the top k features from the learning model. Some examples of the filter FS techniques are the symmetric uncertainty [3], Relief, Fisher, Mutual Information [4], recursive feature elimination, minimum redundancy maximum relevance, Distributed FS [5]. The classifier-dependent techniques are the time taking approaches as it needs a few learning algorithms to select the best features, which reduces the accuracy and the performance [6].

We have assessed the performance of the filter FS techniques using the machine learning techniques like linear regression, decision tree, random forest, XGBoost, lasso regression, and clustering techniques. Clustering techniques help to maintain the stability of a filter FS technique to perform similarity while selecting the subset of features with the same cardinality. In this paper, we selected the best features by measuring the stability of a method as its kurtosis of effectiveness across the features in the dataset, and we compare the stability and the performance of the seven filter FS methods. In this paper, we are using the target feature to remove the irrelevant features.

The summary of the proposed work is:

- We performed filter FS techniques using the symmetric uncertainty and found the minimal subset of features.
- We applied the threshold and constructed the correlation coefficient matrix to measure the dependency between the features.
- We have selected the best features by congregating the features into clusters.
- We reduced the features and hence, the overfitting that improves the prediction accuracy.

The structure of the rest of the paper is as follows. Section 2 reviews the antecedents of the FS techniques in various fields and discusses some existing FS techniques. Section 3 discusses the novel FS algorithm that helps find the minimal subset of features to predict the air quality index. Section 4 describes the experimental results, analyses the proposed algorithm's behavior, and presents the obtained results. Section 5 concludes the paper.

2 Related work

FSs affect the model construction, and the predefined criterion evaluates the optimization of the feature subset. FS approaches consist of selection, evaluation stopping criterion, and validation. The filter FS measures the relevance between the features, different metrics used for this are correlation-based FS algorithms [7], distance-based FS algorithms [8], statistics-based FS algorithm [9], information theory-based FS algorithms [10]. The FS algorithms are divided into linear [11] and non-linear FS algorithms [12]. Fran et al. [13] proposed the conditional mutual information FS technique that is weakly dependent. Bennasar et al. [14] proposed the joint mutual information, which extends the conditional mutual information, and used the maximum, minimum criteria. Zeng et al. [15] proposed interaction weight FS, which dynamically influences the mutual information of the features and the class labels. Hu et al. [16] proposed the dynamic relevance and joint mutual information to remove the redundant features. Kolli et al. [17] proposed a granular feature multi-variant clustering-based genetic algorithm for feature subset selection. This technique uses the granularity of neighborhood-based rough sets and the fitness values as the threshold to subset features. Sai Prasad et al. [18] proposed a novel left-to-right and right-to-left framework to reduce the features and generate a finite number of unique features.

3 Methodology

In this section, we proposed a novel filter FS algorithm, which combines the symmetric uncertainty and the kurtosis to select the features and obtain the low redundant features. Figure 1 shows the proposed technique diagram.

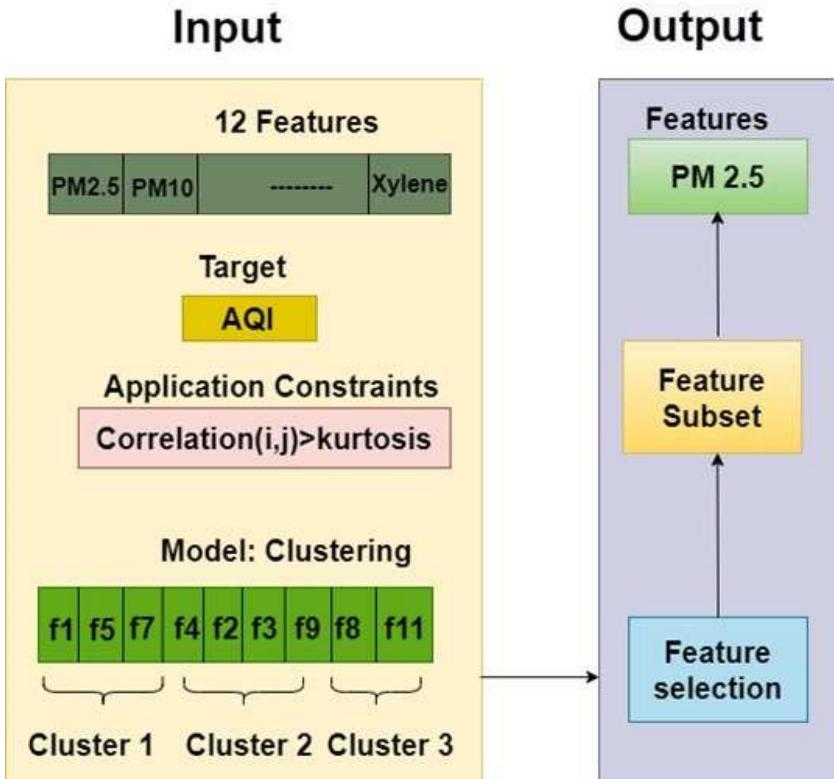


Figure 1. Proposed Block Diagram

In algorithm 1, first, we calculate the symmetric uncertainty values in step-1. In step-2, we calculated the kurtosis of all these symmetric uncertainty values and used this kurtosis value as the threshold. In step-3, we generated the correlation coefficient matrix of the original dataset, and in step-4, we constructed the binary matrix by applying the threshold in the correlation matrix. After generating the binary matrix, step-5 counted the values equal to 1 and stored these values in the $t + 1$ column in the binary matrix. Step-6 generates the clusters based on the binary $t + 1$ columns values, i.e., cluster all the similar count areas in the single cluster. In step-7, we remove the redundant features and maintain the relevant ones as the final subset.

Algorithm 1.

Input: Data set D , feature set $F = f_1, \dots, f_k$

Output: Selected feature set S

1. *Calculate the symmetric uncertainty of each feature and arrange all the features in descending order based on the symmetric uncertainty values.*
2. *Choose the kurtosis value as the threshold.*
3. *Find the correlation coefficient symmetric matrix for the original dataset.*
4. *Apply the threshold value of the generated correlation coefficient matrix.*
 - (a) *If the individual value of the correlation coefficient matrix is greater than the threshold, place the value equal to 1; otherwise, the value equal to 0.*
 - (b) *Repeat the procedure for all the features in the correlation matrix.*
5. *Calculate the total number of ones in each row.*
6. *Combine all the features and form the clusters that have the same weights.*
7. *Choose the highest symmetric uncertainty feature in each cluster.*

4 Results

This paper uses the data collected from 270 monitoring stations from the Indian government website CPCB (Central Pollution Control Board); these stations automatically collect hourly air quality 24 hours per day. The data are open to the public. We collected significant air pollutants, i.e., PM2.5, PM10, CO, NO2, SO2, O3 data, from January 1, 2015, to September 1, 2019 [26]. We used Keras deep learning application programming interface with TensorFlow back end, and we implemented an improved algorithm using the IDE Anaconda.

We applied the proposed kurtosis-based FS(KBFS) algorithm to the collected dataset and first calculated the symmetric uncertainty on the collected air pollution dataset. We considered 12 features (PM 2.5, PM10, CO, NO, NO2, NOx, SO2, O3, NH3, Benzene, Toluene, Xylene) as the input and AQI as the target feature. Figure 2 shows the symmetric uncertainty values of the input features concerning the target feature.

	SU	Feature
0	0.420020	PM2.5
3	0.392949	NO2
8	0.392626	O3
4	0.361884	NOx
2	0.341796	NO
7	0.306341	SO2
1	0.192805	PM10
5	0.132793	NH3
6	0.012845	CO
10	0.004248	Toluene
9	-0.057432	Benzene
11	-0.640781	Xylene

Figure 2. Sorted symmetric uncertainty

After calculating the symmetric uncertainty, we generated the correlation coefficient matrix on the original dataset. Figure 3 shows the correlation coefficient matrix for the air pollution dataset.

After generating the correlation coefficient matrix, calculate the kurtosis for the symmetric uncertainty values according to Figure 2. Set this value as the threshold value. Apply this threshold value to the correlation coefficient matrix. If the coefficient matrix value is greater than the threshold, then set value one; otherwise, set value zero. Figure 4 shows the generated binary matrix.

Count the number of ones in the binary matrix and record those

	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	Toluene	Xylene
PM2.5	1.0000	0.3684	0.4692	0.4332	0.3822	0.1565	0.1128	0.1642	0.2876	0.0458	0.1678	0.1081
PM10	0.3684	1.0000	0.4045	0.3302	0.3822	0.2259	-0.0504	0.1486	0.2810	0.0422	0.1087	0.0167
NO	0.4692	0.4045	1.0000	0.4989	0.7375	0.1710	0.2335	0.2159	0.1210	0.0569	0.1703	0.1006
NO2	0.4332	0.3302	0.4989	1.0000	0.5924	0.1696	0.3692	0.4323	0.3932	0.0643	0.3300	0.2198
NOx	0.3822	0.3822	0.7375	0.5924	1.0000	0.1574	0.2447	0.2182	0.1652	0.0677	0.2088	0.1130
NH3	0.1565	0.2259	0.1710	0.1696	0.1574	1.0000	-0.0736	-0.0578	0.1507	0.0255	0.0196	-0.0528
CO	0.1128	-0.0504	0.2335	0.3692	0.2447	-0.0736	1.0000	0.4780	0.0718	0.0695	0.2908	0.1950
SO2	0.1642	0.1486	0.2159	0.4323	0.2182	-0.0578	0.4780	1.0000	0.2410	0.0494	0.2849	0.2669
O3	0.2876	0.2810	0.1210	0.3932	0.1652	0.1507	0.0718	0.2410	1.0000	0.0507	0.1756	0.1088
Benzene	0.0458	0.0422	0.0569	0.0643	0.0677	0.0255	0.0695	0.0494	0.0507	1.0000	0.6906	0.0967
Toluene	0.1678	0.1087	0.1703	0.3300	0.2088	0.0196	0.2908	0.2849	0.1756	0.6906	1.0000	0.3091
Xylene	0.1081	0.0167	0.1006	0.2198	0.1130	-0.0528	0.1950	0.2669	0.1088	0.0967	0.3091	1.0000

Figure 3. Correlation matrix

	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	Toluene	Xylene
PM2.5	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
PM10	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
NO	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
NO2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
NOx	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
NH3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
CO	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
SO2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
O3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Benzene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Toluene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Xylene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

Figure 4. Binary matrix

counts in the t+1 column in the binary matrix. Figure 5 shows the count for the binary matrix for each input feature.

	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	Toluene	Xylene	count
PM2.5	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
PM10	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
NO	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
NO2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
NOx	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
NH3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
CO	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
SO2	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
O3	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
Benzene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
Toluene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12
Xylene	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	12

Figure 5. Count in the binary matrix using the kurtosis as the threshold

Generate the clusters that have similar counts and select the relevant feature from each cluster. We sorted all the features in the cluster as per its symmetric uncertainty, and then we selected the best symmetric uncertainty as to the relevant feature. Figure 6 shows the final subset of features.

	Feature	SU	count
0	PM2.5	0.42002	12

Figure 6. Subset of features

Table 1 tabulates different statistical tests used to find the threshold we applied to the correlation matrix. These experimental results showed that the final subset of features is only one for kurtosis, so it is the best statistical test measure to set as the threshold.

Table 1. Different statistical tests used for threshold in symmetric uncertainty

Statistical Techniques	Tech-	Total Features	Selected Features
Mean		12	8
Median		12	6
Mode		12	6
Variance		12	6
Standard deviation		12	6
kurtosis		12	1

Table 2 tabulates the comparison of the different filter FS techniques. From the results, we observed that selected features using the symmetric uncertainty are smaller than the remaining filter FS techniques.

Table 2. Comparison of different filter FS techniques

FS Techniques	Total Features	Selected Features
Anova	12	6
Correlation Feature selector	12	5
mRmR	12	6
Fisher	12	5
Mutual Information	12	2
Information Gain	12	4
Gain ratio	12	3
ReliefF	12	3
Variance	12	5
Symmetric Uncertainty	12	1

We reduced the number of features using the filter FS methods, keeping the relevant features that help predict the accurate air quality index. In comparison, the proposed technique selects the most relevant features, and the remaining FS techniques keep an average of 42% of the original features, but the kurtosis-based FS technique reduces the features up to 83.33%, whereas the remaining filter feature selection reduction rate ranges from 40% to 50%. Figure 7 shows the comparison of reduced rates of features using various FS techniques.

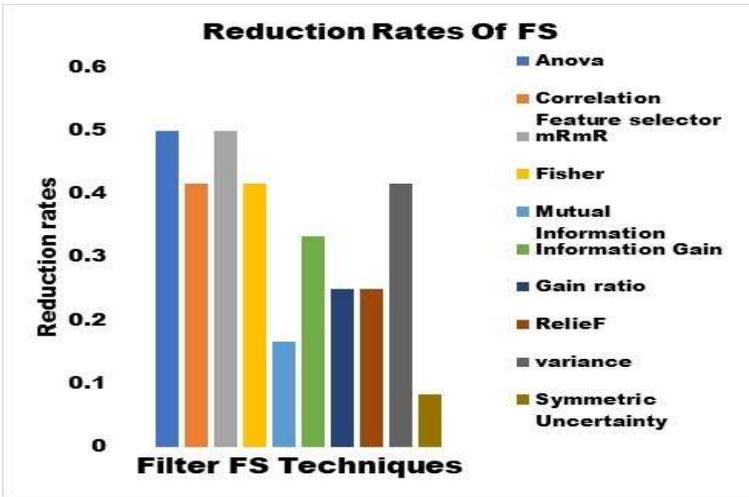


Figure 7. Comparison of reduction rates of features using various FS techniques

We analyzed the air quality dataset using machine learning techniques. We also analyzed the complete air pollution dataset, and the features were selected using the kurtosis-based FS by evaluating the performance metrics like correlation coefficient, root-mean-square error, and accuracy. We implemented the proposed kurtosis-based FS algorithm on different machine learning algorithms and tabulated it in Table 3.

We discussed the performance comparison of the machine learning classifiers after performing the FS. Table 4 compares the different FS techniques using the different classifiers. This paper uses the classifiers Random Forest(RF), Linear Regression(LR), and Principal Component

Table 3. Comparison of different machine learning algorithms with proposed technique

Classifiers	Without FS				With FS			
	r	R2	RMSE	ACC	r	R2	RMSE	ACC
Decision Tree	0.847	0.717	37.07	76.95	0.894	0.826	29.53	91.76
Linear Regression	0.827	0.736	35.45	74.53	0.874	0.825	31.42	89.34
Lasso regression	0.814	0.717	37.83	78.58	0.872	0.838	29.32	89.53
Random Forest	0.864	0.767	36.34	78.56	0.914	0.848	21.57	93.78
XGBoost	0.826	0.762	33.56	72.45	0.893	0.852	27.45	89.74
Support Vector machine	0.823	0.736	33.53	72.54	0.864	0.857	25.53	88.86

Analysis(PCA).

From the results, we observed that processing time is less for the proposed technique. We observed that accuracy significantly improved by applying the proposed technique and accomplished better performances for the symmetric uncertainty for the random forest classifier, and the processing time is 09ms.

We assessed the correlation between the major air pollutants PM2.5, CO, NO2, O3, SO2, and the air quality index and finally explored their relationship. Figure 8 shows the correlation between the major air pollutants and the air quality index.

This research finds the minimal subset of features that helps predict the accurate air quality index. From the observations, we found that PM2.5 is a relevant feature to predict the air quality index.

5 Conclusion

The objective of the FS is to select a minimal subset of relevant features. We proposed a kurtosis-based FS algorithm to reduce the dimensionality of the air pollution data by selecting the best features, which enhances the prediction performance. This paper discusses different filter FS techniques and various statistical tests to finalize the threshold to find the best fit subset of features. The authors performed the comparison for different filter FS techniques, the results showed that the proposed kurtosis-based FS algorithm improves the prediction performance.

Table 4. Comparison of different FS techniques using different classifiers

Classifier	Measure	r	R2	RSME	ACC	PT(ms)
RF	Chi-square [19]	0.825	0.701	36.54	85.78	13
	Relief [20]	0.802	0.703	37.73	84.62	14
	MultiSurf [21]	0.826	0.708	36.86	85.34	13
	Ensemble FS [22]	0.904	0.748	31.57	83.78	11
	Anova [23]	0.809	0.717	38.45	84.03	14
	M-Cluster FS [24]	0.895	0.826	22.45	91.34	10
	SU-MLP [25]	0.897	0.829	22.13	92.31	10
	Proposed	0.914	0.848	21.57	93.78	09
	LR	Chi-square [19]	0.818	0.716	39.69	85.27
Relief [20]		0.817	0.701	41.96	84.93	13
MultiSurf [21]		0.814	0.701	40.51	85.34	12
Ensemble FS [22]		0.910	0.741	36.54	83.53	11
Anova [23]		0.813	0.704	39.34	85.45	14
M-Cluster FS [24]		0.862	0.817	32.68	87.68	13
SU-MLP [25]		0.869	0.821	31.79	88.95	13
Proposed		0.874	0.825	31.42	89.34	12
PCA		Chi-square [19]	0.834	0.705	25.52	86.19
	Relief [20]	0.827	0.719	26.34	85.82	12
	MultiSurf [21]	0.829	0.715	25.64	86.45	13
	Ensemble FS [22]	0.918	0.749	22.49	86.57	12
	Anova [23]	0.827	0.721	27.54	85.73	12
	M-Cluster FS [24]	0.884	0.815	21.53	90.27	12
	SU-MLP [25]	0.889	0.821	20.84	90.95	12
	Proposed	0.894	0.826	20.53	91.76	11

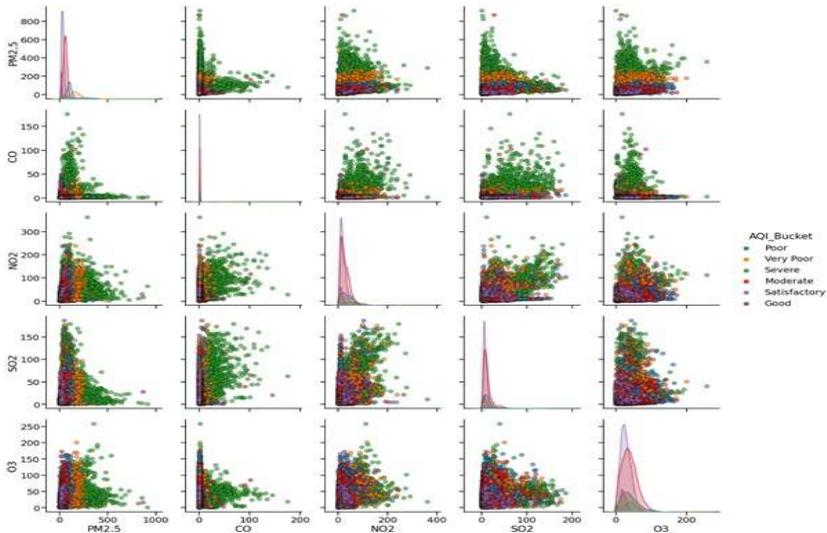


Figure 8. Correlation between the major air pollutants and the air quality index

References

- [1] X. Wang, B. Guo, Y. Shen, C. Zhou, and X. Duan, “Input feature selection method based on feature set equivalence and mutual information gain maximization,” *IEEE Access*, vol. 7, no. 1, pp. 151525–151538, 2019. DOI: 10.1109/ACCESS.2019.2948095.
- [2] F. Macedo, M. R. Oliveira, A. Pacheco, and R. Valadas, “Theoretical foundations of forward feature selection methods based on mutual information,” *Neurocomputing*, vol. 325, pp. 67–89, 2019.
- [3] S. P. Potharaju and M. Sreedevi, “A novel cluster of quarter feature selection based on symmetrical uncertainty,” *Gazi University Journal of Science.*, vol. 31, no. 2, pp. 456–470, 2018.
- [4] E. Hancer, B. Xue, and M. Zhang, “Differential evolution for filter feature selection based on information theory and feature ranking,” *Knowledge-Based Systems*, vol. 140, pp. 103–119, 2018.
- [5] S. P. Potharaju and M. Sreedevi, “Distributed feature selection (DFS) strategy for microarray gene expression data to improve

- the classification performance,” *Clinical Epidemiology and Global Health*, vol. 7, no. 2, pp. 171–176, 2019.
- [6] G. Chandrashekar and F. Sahin, “A survey on feature selection methods,” *Computers & Electrical Engineering*, vol. 40, no. 1, pp. 16–28, 2014.
- [7] J. Xie, M. Wang, and Q. Hu, “The differentially expressed gene selection algorithms for unbalanced gene datasets by maximize the area under ROC,” *Journal of Shaanxi Normal University (Natural Science Edition)*, vol. 1, no. 1, pp. 01–11, 2017.
- [8] Y. Sun and J. Li, “Iterative RELIEF for feature weighting,” in *Proceedings of the 23rd international conference on Machine learning*, vol. 1, no. 1, pp. 913–920, 2006.
- [9] J. Dai and Q. Xu, “Attribute selection based on information gain ratio in fuzzy rough set theory with application to tumor classification,” *Applied Soft Computing*, vol. 13, no. 1, pp. 211–221, 2013.
- [10] S. Sharmin, M. Shoyaib, A. A. Ali, M. A. H. Khan, and O. Chae, “Simultaneous feature selection and discretization based on mutual information,” *Pattern Recognition*, vol. 19, no. 29, pp. 162–174, 2019.
- [11] J. Che, Y. Yang, L. Li, X. Bai, S. Zhang, and C. Deng, “Maximum relevance minimum common redundancy feature selection for non-linear data,” *Information Sciences*, vol. 409, pp.68–86, 2017.
- [12] P. Zhang, W. Gao, and G. Liu, “Feature selection considering weighted relevancy,” *Applied Intelligence*, vol. 48, no. 12, pp. 4615–4625, 2018.
- [13] F. Fleuret, “Fast binary feature selection with conditional mutual information,” *Journal of Machine learning research*, vol. 5, no. 9, pp. 1–11, 2014.
- [14] M. Bannasar, Y. Hicks, and R. Setchi, “Feature selection using joint mutual information maximisation,” *Expert Systems with Applications*, vol. 42, no. 22, pp. 8520–8532, 2015.
- [15] Z. Zeng, H. Zhang, R. Zhang, and C. Yin, “A novel feature selection method considering feature interaction,” *Pattern Recognition*, vol. 48, no. 8, pp. 2656–2666, 2015.

- [16] Hu, Liang and Gao, Wanfu and Zhao, Kuo and Zhang, Ping and Wang, Feng. “Feature selection considering two types of feature relevancy and feature interdependency,” *Expert Systems with Applications*, vol. 93, pp. 423–434, 2018.
- [17] Kolli, Srinivas, M.Sreedevi, “A Novel Granularity Optimal Feature Selection based on Multi-Variant Clustering for High Dimensional Data,” *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, vol.12, no. 3, pp. 5051–5062, 2021.
- [18] Potharaju, Sai Prasad and Sreedevi, M, “A novel LtR and RtL framework for subset feature selection (reduction) for improving the classification accuracy,” *Advanced Computing and Intelligent Engineering*, pp. 215–224, 2019.
- [19] Jin, Xin and Xu, Anbang and Bie, Rongfang and Guo, Ping, “Machine learning techniques and chi-square feature selection for cancer classification using SAGE gene expression profiles,” *International workshop on data mining for biomedical applications*, pp. 106–115, 2016.
- [20] Mesin, Luca and Orione, Fiammetta and Taormina, Riccardo and Pasero, Eros, M, “A feature selection method for air quality forecasting,” *International Conference on Artificial Neural Networks*, pp. 489–494, 2010.
- [21] Urbanowicz, Ryan J and Meeke, Melissa and La Cava, William and Olson, Randal S and Moore, Jason H, “Relief-based feature selection: Introduction and review,” *Journal of biomedical informatics*, vol. 85, pp. 189–203, 2018.
- [22] Chen, Chih-Wen and Tsai, Yi-Hong and Chang, Fang-Rong and Lin, Wei-Chao, “Ensemble feature selection in medical datasets: Combining filter, wrapper, and embedded feature selection results,” *Expert Systems*, vol.37, no. 15, pp. 1–15, 2020.
- [23] Ding, Hui and Feng, Peng-Mian and Chen, Wei and Lin, Hao, “Identification of bacteriophage virion proteins by the ANOVA feature selection and analysis,” *Molecular BioSystems*, vol.10, no. 8, pp. 2229–2235, 2014.
- [24] Potharaju, Sai Prasad and Sreedevi, M, “A Novel M-Cluster of Feature Selection Approach Based on Symmetrical Uncertainty for

- Increasing Classification Accuracy of Medical Datasets,” *Journal of Engineering Science & Technology Review*, vol.10, no. 6, pp. 1-8, 2019.
- [25] Potharaju, Sai Prasad and Sreedevi, M and Amiripalli, Shanmuk Srinivas, M, “An ensemble feature selection framework of sonar targets using symmetrical uncertainty and multi-layer perceptron (su-mlp),” *Cognitive Informatics and Soft Computing*, pp. 247–256, 2019.
- [26] Ministry of Environment Forest and Climate CHange, Government of India, “Central Pollution Control Board,” <https://cpcb.nic.in/>, accessed December 7, 2020.

Usharani Bhimavarapu, M. Sreedevi

Received January 21, 2022

Accepted June 23, 2022

Usharani Bhimavarapu
Department of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram, Andhra Pradesh, India.
E-mail: ushareddy@kluniversity.in

M. Sreedevi
Department of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram, Andhra Pradesh, India.
E-mail: msreedevi27@kluniversity.in

Wiener Index of Some Brooms

Julian D. Allagan

Abstract

In the field of chemical graph theory, a Wiener (topological) index is a type of a molecular descriptor that is calculated based on the molecular graph of alkanes. It gives the sum of geodesic distances (or shortest paths) between all pairs of vertices of the graph. We found and prove the Wiener indices of some Brooms, which are Caterpillars, giving several unknown sequences that are now added to the collection of the largest Online Encyclopedia of Integer Sequences.

Keywords: Wiener Index, Brooms, Sequences.

1 Introduction

Suppose G is a simple graph and $v \in V(G)$. The distance between two vertices $u, v \in V(G)$, often denoted by $d_G(u, v)$, is the length (number of edges) of their shortest path in G ; this is also known as a *geodesic distance*. The *eccentricity* of a vertex v , written as $\epsilon(v)$, is the maximum of the distance between v and any other vertex $u \in V$, i.e., $\epsilon(v) = \max_{u \in V} \{d_G(v, u)\}$. Further, the *diameter* d of a graph is the maximum eccentricity of any vertex in the graph, i.e., $d = \max_{v \in V} \epsilon(v)$. These parameters are often useful in classifying acyclic (tree-like) graphs. A topological index, a numerical value, is often used to describe chemical structures in the areas of chemical graph theory, molecular topology, and mathematical chemistry. One of the well-studied indices is the Wiener topological index (W), introduced in 1947 by Harry Wiener [12]. In graph theory, the *Wiener index* of a graph G , denoted by $W(G)$, is the sum of the distances between all unordered pairs of vertices of G . This index is used to describe and even predict

several physical and chemical properties of molecules such as density, viscosity and velocity. Details of some of these reviews can be found in [4], [6], [8], [9], [11], for instance.

Throughout this article, we denote the n^{th} triangular number by $T_n = \frac{n(n+1)}{2} = 1 + 2 + 3 + 4 + 5 + 6 + \dots + n$. Moreover, we denote the k^{th} tetrahedral number by $T_n^k = \sum_k T_n$, the sum of the first k n^{th} triangular numbers.

Suppose $P_n := v_1 - v_2 - \dots - v_{n-1} - v_n$, denotes a path on $n \geq 3$ vertices. By sequentially connecting k leaves to some vertex v_i , with $2 \leq i \leq n-1$, we obtain a *Caterpillar*. We denote a Caterpillar on $n+k$ vertices by P_n^k , where P_n is referred to as *stem* or *backbone*; observe that the diameter of P_n^k is the length of P_n . If every internal vertex is adjacent to at least one of the k new pendant vertices, then P_n^k is said to be *complete*. Caterpillars have been used in chemical graph theory to represent the structure of benzenoid hydrocarbon molecules, e.g., [3] and [5]. Due to their importance, we present the formulae of several Wiener values for some Caterpillars. In particular, we found that the Wiener values, as sequences, of these Caterpillars which contain exactly one vertex of degree greater than 2 do not currently exist in the largest Online Encyclopedia of Integer Sequences (OEIS) [10]. We hope to submit these values and their formulae for the record.

It is easy to see that, for any complete graph K_n , $W(K_n) = \binom{n}{2}$. Further, it is well-known that the Wiener indices of a star graph S_n and a path graph P_n are $(n-1)^2$ and $\frac{(n-1)n(n+1)}{6}$, respectively. See [1] and [5], for instance. Further, it is shown [4] that, for any tree T on n vertices, $W(S_n) < W(T) < W(P_n)$.

As an example for finding Wiener indices, we present a proof for the Wiener index of a path, after the next proposition.

Proposition 1. $\sum_{j=2}^n \binom{j}{2} = \binom{n+1}{3}$ holds for all $n \geq 2$

Proof. The case when $n = 2$ is trivial. Let's assume for all $k \geq 2$,

$\sum_{j=2}^k \binom{j}{2} = \binom{k+1}{3}$. It follows that

$$\begin{aligned} \sum_{j=2}^{k+1} \binom{j}{2} &= \sum_{j=2}^k \binom{j}{2} + \binom{k+1}{2} \\ &= \binom{k+1}{3} + \binom{k+1}{2} \\ &= \binom{k+2}{3}. \end{aligned}$$

Hence the result by induction. □

Corollary 1. *The Wiener index of a path P_n is $W(P_n) = \binom{n+1}{3}$, $n \geq 2$.*

Proof. Suppose the vertices of the path are v_1, v_2, \dots, v_n . We proceed to add the distances between v_i, v_j , for each $i \neq j$. As such, we compute $\sum_{\substack{j>i \\ j \leq n}} d(v_i, v_j)$ which is equal to $\binom{n-i+1}{2}$, for each i , with $1 \leq i \leq n$.
Now,

$$\begin{aligned} W(P_n) &= \sum_{i=1}^{n-1} \sum_{j>i}^n d(v_i, v_j) \\ &= \sum_{i=1}^{n-1} \binom{n-i+1}{2} \\ &= \sum_{j=2}^n \binom{j}{2}, \end{aligned}$$

for all $n \geq 2$. The result follows from Proposition 1.

Alternatively, given the n^{th} triangular number T_n , we have

$$\begin{aligned}
 W(P_n) &= \sum_{j=1}^{n-1} T_j = \sum_{j=1}^{n-1} \frac{j(j+1)}{2} \\
 &= \frac{1}{2} \left(\sum_{i=1}^{n-1} j^2 + \sum_{i=1}^{n-1} j \right) \\
 &= \frac{1}{2} \left(\frac{n(n-1)[2(n-1)+1]}{6} + \frac{n(n-1)}{2} \right) \\
 &= \frac{(n-1)n(n+1)}{6} \\
 &= \binom{n+1}{3}.
 \end{aligned}$$

□

Remark 1.

Recall that, the n^{th} *rising factorial* and the n^{th} *falling factorial* denoted respectively by $x^{\overline{n}}$ and $x^{\underline{n}}$, are $x^{\overline{n}} = x(x+1)(x+2) \cdots (x+n-1) = \prod_{k=1}^n (x+k-1) = \prod_{k=0}^{n-1} (x+k)$ and $x^{\underline{n}} = x(x-1)(x-2) \cdots (x-n+1) = \prod_{k=1}^n (x-k+1) = \prod_{k=0}^{n-1} (x-k)$. Although the previous result (and upcoming ones) can be written in either format, i.e., $W(P_n) = \frac{(n-1)^{\overline{3}}}{3!} = \frac{(n+1)^{\underline{3}}}{3!}$, it is beyond the interest of this article.

2 Wiener index of Comb graphs

Suppose $P_n := v_1 - v_2 - \dots - v_{n-1} - v_n$ denotes a path on $n \geq 3$ vertices. By sequentially connecting a single leaf to each vertex v_i , with $2 \leq i \leq n-1$, we obtain a *Complete Caterpillar* P_n^{n-2} which is commonly known as a *Comb* or a *Centipede*. Figure 2 is an example.

For the next result, we define the following: Suppose $G = (V, E)$ denotes a graph with an ordered list of vertices (v_1, v_2, \dots, v_n) . We

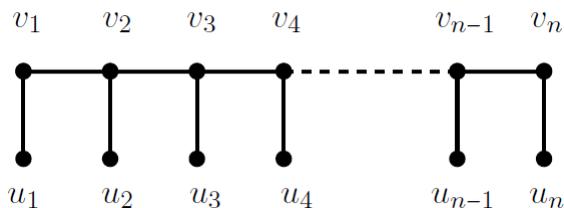


Figure 1. Comb graph on $2n$ vertices

denote and define the s^{th} partial Wiener (index) of G by $W^s(G)$ and $W^s(G) = \sum_{j>s} d(v_s, v_j)$, respectively. Thus, the 1^{st} Partial Wiener (index) of G is the sum of the distances between v_1 and any other vertex $v_j \in V, j > 1$. The special case when $G = K_2, W^1(G) = 1 = W(G)$.

Lemma 1. *Suppose G is any graph with an ordered list of vertices (v_1, v_2, \dots, v_n) . It follows that $W(G) = \sum_{k=1}^{n-1} \sum_{j>s} d(v_s, v_j)$.*

Proof. From the definitions of s^{th} partial Wiener and the (full) Wiener index of G , it follows that $W(G) = \sum_{s=1}^{n-1} W^s(G)$, giving the result. \square

Theorem 2. *Suppose G^n denotes a Comb on $2n$ vertices. Then $W(G^n) = \frac{n(2n^2 + 6n - 5)}{3}$ for all $n \geq 1$.*

Proof. Let $G^1 := u_1 - v_1$, a path on 2 vertices. We add two pendant vertices u_2, v_2 , such that u_2 is adjacent to v_2 and v_2 is adjacent to v_1 . The resulting graph is a Comb, denoted by G^2 , which is isomorphic to P_4 . Thus, $W^1(G^1) = d(u_1, v_1) = T_1$ and $W^2(G^1) = T_0$ since $W^2(G^1) = W^1(G^1 \setminus \{u_1\})$. Further, $W^1(G^2) = \sum_{w \in V(G^2)} d(u_2, w) = T_3$ and $W^2(G^2) = \sum_{w \in V(G^2)} d(u_2, w) = T_2$. So, when $G = P_4$, by definition

of partial Wieners,

$$\begin{aligned} W(G) &= W^1(G^1) + W^2(G^1) + W^1(G^2) + W^2(G^2) \\ &= T_0 + T_1 + T_2 + T_3 \\ &= 10. \end{aligned}$$

Iteratively, for each $k \geq 1$, we form G^k from a previously formed Comb G^{k-1} , by adding the pair of vertices (u_k, v_k) such that u_k is adjacent to v_k and v_k is adjacent to $v_{k-1} \in V(G^{k-1})$. Thus, the vertices of G^k can be seen as the ordered list $(u_1, v_1, u_2, v_2, \dots, u_k, v_k)$. With each such pair $(u_k, v_k) \notin V(G^{k-1})$ we compute and add the first and second partial Wieners of G^k , $k \geq 3$. So, given u_k , $W^1(G^k) = \sum_{w \in V(G^k)} d(u_k, w)$.

Because $d(u_k, u_{k-1}) = T_3$ and $d(u_k, u_1) = T_{k+1}$, for all $k \geq 2$, it follows that

$$\begin{aligned} W^1(G^k) &= T_3 + \sum_{j=1}^{k-2} (T_{3+j} - T_{1+j}) \\ &= (T_{k+1} - T_{k-1}) + (T_k - T_{k-2}) + \dots + (T_5 - T_3) + \\ &\quad + (T_4 - T_2) + T_3 \\ &= T_{k+1} + T_k - T_2, \quad k \geq 3. \end{aligned}$$

Similarly, given (u_k, v_k) , we compute the second partial Wieners of G^k , i.e., $W^2(G^k) = \sum_{w \in V(G^k)} d(v_k, w)$. Because $d(v_k, u_{k-1}) = T_2$ and $d(v_k, u_1) = T_k$, for all $k \geq 3$, it follows that

$$\begin{aligned} W^2(G^k) &= T_2 + \sum_{j=1}^{k-2} (T_{2+j} - T_j) \\ &= (T_k - T_{k-2}) + (T_{k-1} - T_{k-3}) + \dots + (T_4 - T_2) + \\ &\quad + (T_3 - T_1) + T_2 \\ &= T_k + T_{k-1} - T_1, \quad k \geq 3. \end{aligned}$$

Therefore, for all $n \geq 1$,

$$\begin{aligned}
 W(G) &= \sum_{k=1}^n W^1(G^k) + \sum_{k=1}^n W^2(G^k) \\
 &= \sum_{k=1}^2 W^1(G^k) + \sum_{k=1}^2 W^2(G^k) + \sum_{k=3}^n W^1(G^k) + \sum_{k=3}^n W^2(G^k) \\
 &= T_1 + T_2 + T_3 + \sum_{k=3}^n W^1(G^k) + \sum_{k=3}^n W^2(G^k) \\
 &= T_1 + T_2 + T_3 + \sum_{k=3}^n (T_{k+1} + 2T_k + T_{k-1} - T_2 - T_1) \\
 &= T_1 + T_2 + T_3 + \sum_{k=3}^n (T_{k+1} + 2T_k + T_{k-1}) - \sum_{k=3}^n (T_2 + T_1).
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 W(G) &= T_1 + T_2 + T_3 + \sum_{k=1}^{n+1} T_k + 2 \sum_{k=1}^n T_k \\
 &+ \sum_{k=1}^{n-1} T_k - \left(\sum_{k=1}^3 T_k + 2 \sum_{k=1}^2 T_k + T_1 \right) - (n-2)(T_2 + T_1) \\
 &= 3T_n + T_{n+1} + 4 \sum_{k=1}^{n-1} T_k - (4n+1) \\
 &= \frac{3n(n+1)}{2} + \frac{(n+1)(n+2)}{2} + 2n(n-1) - 4n - 1 \\
 &= \frac{2n^3}{3} + 2n^2 - \frac{5n}{3},
 \end{aligned}$$

giving the result for all $n \geq 1$. □

Here, in Table 2, we present the first ten values of the Wiener of Combs. We note that Emeric Deutsch had submitted (in 2011) this formula to OEIS as **A192023** [10] and yet, we have no record of the proof of the result.

Table 1. The first ten values of the Wiener of a Comb graph on $2n$ vertices

n	1	2	3	4	5	6	7	8
$\frac{2n^3}{3} + 2n^2 - \frac{5n}{3}$	1	10	31	68	125	206	315	456

3 Wiener Index of Brooms

A Caterpillar that is obtained by adding $k \geq 1$ pendant vertices to the first (or last) internal vertex of P_n is called a *Broom*. We denote it as B_n^k . The graphs generated in the special cases when $k = 1$ and $k = 2$ are called, respectively, *Sling* and *Tridon*. Figure 2 shows a Tridon.

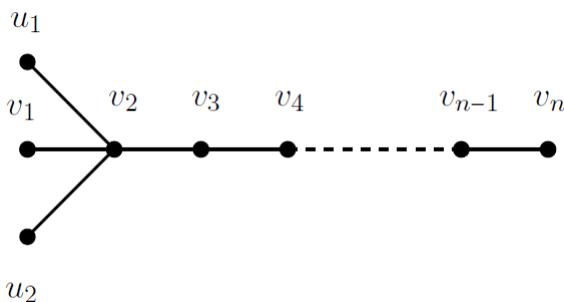


Figure 2. A Tridon B_n^2 on $n + 2$ vertices

Theorem 3. Suppose B_n^k denotes a Broom on $n + k$ vertices. Then $W(B_n^k) = 2T_k + kT_{n-1} + W(P_n)$ for all $n \geq 3, k \geq 1$.

Proof. Let T_n denote the n^{th} triangular number and $u_i, i = 1 \dots, k$, the pendant vertices. Consider a main path P_n , for $n \geq 3$, giving $W(P_n)$.

With each additional pendant vertex u_i connected to $v_2 \in P_n$, we have $d(u_i, v) = T_{n-1}$, for each $v \in P_n$ and $v \neq v_1$. This gives $\sum_{j=1}^k T_{n-1}$ for each $u_i, i = 1, \dots, k$. Finally $d(u_i, v_1) = 2 = d(u_i, u_j)$ with $i \neq j$,

giving $\sum_{j=1}^k 2j$. Together, we have

$$\begin{aligned} W(B_n^k) &= \sum_{j=1}^k 2j + \sum_{j=1}^k T_{n-1} + W(P_n) \\ &= 2 \sum_{j=1}^k j + \sum_{j=1}^k T_{n-1} + W(P_n) \\ &= 2T_k + kT_{n-1} + W(P_n), \end{aligned}$$

giving the result. □

Corollary 2. *If B_n^n denotes a Broom on $2n$ vertices of which $n+2$ are pendant, then $W(B_n^n) = \frac{2n^3}{3} + \frac{n^2}{2} + \frac{5n}{6}$, for all $n \geq 3$ vertices.*

Proof. From Theorem 3 when $k = n$, there are exactly $n + 2$ pendant vertices and we have

$$\begin{aligned} W(B_n^1) &= 2T_n + nT_{n-1} + W(P_n) \\ &= n(n+1) + \frac{n^2(n-1)}{2} + \frac{(n+1)n(n-1)}{6} \\ &= \frac{2n^3}{3} + \frac{n^2}{2} + \frac{5n}{6}. \end{aligned}$$

□

This formula and several of its values are submitted and they are approved in OEIS [10] as **A349416**. In Table 3, we list the first ten values.

Corollary 3. *If B_n^1 denotes a Sling graph, then $W(B_n^1) = \frac{n^3}{6} + \frac{n^2}{2} - \frac{2n}{3} + 2$, for all $n \geq 3$ vertices.*

Table 2. The first eight values of the Wiener of a Broom graph on $2n$ vertices (of which $n + 2$ are pendant)

n	3	4	5	6	7	8	9	10
$\frac{2n^3}{3} + \frac{n^2}{2} + \frac{5n}{6}$	25	54	100	167	259	380	534	725

Proof. By definition, a Sling is a Broom on $k = 1$ pendant vertex. So, when $k = 1$, the result in Theorem 3 becomes

$$\begin{aligned}
 W(B_n^1) &= 2T_1 + T_{n-1} + W(P_n) \\
 &= 2 + \frac{(n-1)n}{2} + W(P_n) \\
 &= 2 + \frac{n(n-1)}{2} + \frac{(n+1)n(n-1)}{6} \\
 &= \frac{n^3}{6} + \frac{n^2}{2} - \frac{2n}{3} + 2.
 \end{aligned}$$

□

This formula is now approved in OIES as **A349417**. We found later that such values are also equivalent to the sequence **A005581**+2 which carries many combinatorics and algebraic meanings. For instance, **A005581** gives the number of inscribable triangles within a $(n + 4)$ -gon sharing with them its vertices but not its sides, according to Lekraj Beedassy [10].

In Table 3, we present the first ten values of the Wiener index of a Sling.

Table 3. The first eight values of the Wiener of a Sling graph on $n + 1$ vertices

n	3	4	5	6	7	8	9	10
$\frac{n^3}{6} + \frac{n^2}{2} - \frac{2n}{3} + 2$	9	18	32	52	79	114	158	212

Corollary 4. *If B_n^2 is a Tridon graph, then $W(B_n^2) = \frac{n^3}{6} + n^2 - \frac{7n}{6} + 6$ for all $n \geq 3$.*

Proof. By definition, we obtain the Wiener value of a Tridon from Theorem 3, when $k = 2$, in which we have

$$\begin{aligned} W(B_n^2) &= 2T_2 + 2T_{n-1} + W(P_n) \\ &= 6 + 2\frac{(n-1)n}{2} + W(P_n) \\ &= 6 + n(n-1) + \frac{(n+1)n(n-1)}{6} \\ &= \frac{n^3}{6} + n^2 - \frac{7n}{6} + 6 \end{aligned}$$

after an expansion. □

This formula and several of its values are submitted and they are approved in OEIS [10] as **A349418** .Table 3 shows the first ten values.

Table 4. The first eight values of the Wiener of a Tridon graph on $n + 2$ vertices

n	3	4	5	6	7	8	9	10
$\frac{n^3}{6} + n^2 - \frac{7n}{6} + 6$	16	28	46	71	104	146	198	261

4 Wiener index of 2-Extended Brooms

Here, we present a generalization of Brooms, by extending the original definition from adding pendant vertices to adding paths. Suppose $P_n := v_1 - v_2 - \dots - v_{n-1} - v_n$ denotes a path on $n \geq 3$ vertices. By sequentially adding some path graph P'_m , on $m \geq 1$ vertices to some v_i , $2 \leq i \leq n - 1$, we obtain an m -Extended Broom which we denote by $B_n^k(m)$. The special case when $m = 1$ is a (regular) Broom, i.e., $B_n^k(1) = B_n^k$. Here, we present the case when $m = 2$.

For the upcoming result, for simplicity, let $G^k = B_n^k(2)$ denote a 2-Extended Broom obtained by adding $k \geq 1$ $P'_k := u_{1k} - u_{2k}$ to $v_2 \in P_n$, for $k \geq 1$. See Figure 3 for the case when $k = 2$.

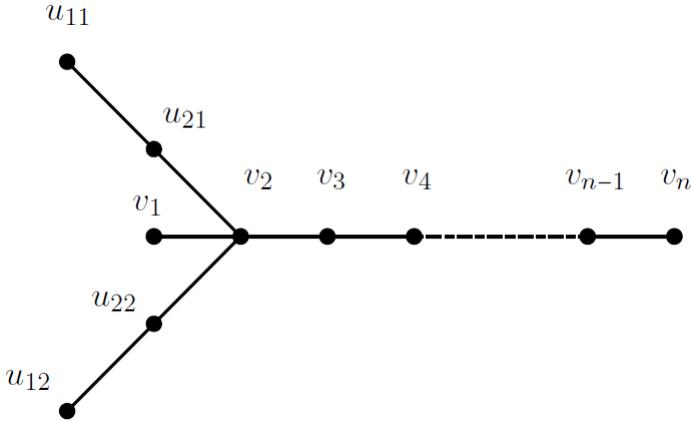


Figure 3. A 2-Extended Broom

Theorem 4. *The Wiener index of a 2-Extended Broom G^k is given by $W(G^k) = \frac{1}{6}n^3 - \frac{1}{6}n + n^2k + 6k^2 - k$ with $n \geq 3$ and $k \geq 1$.*

Proof. Consider the main path, P_n . Now, we sequentially add $P'_k := u_{1k} - u_{2k}$ to $v_2 \in P_n$, for $k \geq 1$. Knowing $W(P_n)$, we proceed to add the values of $d(u_{1k}, v)$ and $d(u_{2k}, v)$, for each $v \in P_n$ and $k \geq 1$.

Observe that $\sum_{x_k} d(u_{1k}, x_k) = T_n$ and $\sum_y d(u_{2k}, y) = T_{n-1}$ for every $x_k \in \{u_{2k}, v_2, v_3, \dots, v_n\}$ and $y \in \{v_2, v_3, \dots, v_n\}$. Further, $d(u_{1k}, v_1) = 3$ and $d(u_{2k}, v_1) = 2$.

Thus, when $k = 1$, we have

$$\begin{aligned} W(G^1) &= W(P_n) + \\ &+ \sum_{x_1} d(u_{11}, x_1) + \sum_y d(u_{21}, y) + d(u_{11}, v_1) + d(u_{21}, v_1) \\ &= W(P_n) + T_n + T_{n-1} + 2(1) + 3(1). \end{aligned}$$

When $k = 2$. We have $d(u_{12}, v_1) = 3 = d(u_{12}, u_{21})$, $d(u_{22}, v_1) = 2 = d(u_{22}, u_{21})$, and $d(u_{12}, u_{11}) = 4$, $d(u_{22}, u_{11}) = 3$. Together, with $\sum_{x_2} d(u_{12}, x_2) + \sum_y d(u_{22}, y)$ for every $x_2 \in \{u_{22}, v_2, v_3, \dots, v_n\}$ and $y \in \{v_2, v_3, \dots, v_n\}$, we have

$$W(G^2) = W(G^1) + \sum_{x_2} d(u_{12}, x_2) + \sum_y d(u_{22}, y) + 2(2) + 3(2) + (3 + 4).$$

Similarly, when $k = 3$, we obtain

$$W(G^3) = W(G^1) + W(G^2) + \sum_{x_3} d(u_{13}, x_3) + \sum_y d(u_{23}, y) + 2(3) + 3(3) + 7(2).$$

Thus, for all $k \geq 1$, we obtain recursively that,

$$\begin{aligned} W(G^k) &= \sum_{i=1}^{k-1} W(G^i) + \sum_{x_k} d(u_{1k}, x_k) + \sum_y d(u_{2k}, y) + 2 \left(\sum_{i=1}^k i \right) \\ &+ 3 \left(\sum_{i=1}^k i \right) + 7 \left(\sum_{i=1}^{k-1} i \right) \\ &= W(P_n) + kT_n + kT_{n-1} + 2T_k + 3T_k + 7T_{k-1} \\ &= W(P_n) + k(T_n + T_{n-1}) + 5T_k + 7T_{k-1}. \end{aligned}$$

Since $W(P_n) = \frac{(n+1)n(n-1)}{6}$ and $T_j = \frac{j(j+1)}{2}$, the result follows after expansion. □

In the next two corollaries, we present two extremal cases; when $k = n$ and when $k = 1$. Both cases follow directly from the previous theorem. We present the first ten values for each case and we point out that neither sequence currently exists in OEIS [10].

Corollary 5. *The Wiener index of a 2-Extended Broom G^n on $3n$ vertices is given by $W(G^n) = \frac{7}{6}n^3 + 6n^2 - \frac{7n}{6}$ with $n \geq 3$.*

Table 5. The first eight values of a 2-Extended Broom G^n on $3n$ vertices

n	3	4	5	6	7	8	9	10
$\frac{7}{6}n^3 + 6n^2 - \frac{7n}{6}$	82	166	290	461	686	972	1326	1755

Table 5 shows some of the values of $W(G^n)$.

Corollary 6. *The Wiener index of a 2-Extended Broom G^1 on $n + 2$ vertices is given by $W(G^1) = \frac{1}{6}n^3 + n^2 - \frac{1}{6}n + 5$ with $n \geq 3$.*

Table 6 shows some of the values of $W(G^1)$.

Table 6. The first eight values of a 2-Extended Broom G^1 on $n + 2$ vertices

n	3	4	5	6	7	8	9	10
$\frac{1}{6}n^3 + n^2 - \frac{1}{6}n + 5$	18	31	50	76	110	153	206	270

References

- [1] D. Bonchev and N. Trinajstić, “Information theory, distance matrix, and molecular branching,” *Journal of Chemical Physics*, vol. 67, no. 10, pp. 4517–4533, 1977.
- [2] E. R. Canfield, R. W. Robinson, and D. H. Rouvray, “Determination of the Wiener molecular branching index for the general tree,” *J. Comput. Chem.*, vol. 6, pp. 598–609, 1985.
- [3] A. Coja-Oghlan, S.O. Krumke, and T. Nierhoff, “Scheduling a server on a caterpillar network—a probabilistic analysis,” in *6th Workshop on Models and Algorithms for Planning and Scheduling Problems*, (France), pp. 48–50, 2003.
- [4] A.A. Dobrynin, R. Entringer, and I. Gutman, “Wiener Index for trees: Theory and Applications,” *Acta Appl. Math.*, vol. 66, pp. 211–249, 2001. DOI: <https://doi.org/10.1023/A:1010767517079>.

- [5] S. El-Basil, "Caterpillar (Gutman) trees in chemical graph theory," *Topics in Current Chemistry*, vol. 153, pp. 273–289, 1990.
- [6] D. Papazova, M. Dimov, and D. Bonchev, "Calculation of gas chromatographic retention indices of isoalkanes based on a topological approach," *Journal of Chromatography*, vol. 188, no. 2, pp. 297–303, 1980.
- [7] D. H. Rouvray and R. B. King, *Topology in Chemistry: Discrete Mathematics of Molecules*, Horwood Publishing Ltd., Chichester, 2002, pp. 89–112.
- [8] D. H. Rouvray and B. C. Crafford, "The dependence of physico-chemical properties on topological factors," *South African Journal of Science*, vol. 72, pp. 47–51, 1976.
- [9] D. H. Rouvray and W. Tatong, "Novel applications of topological indices," *Zeitschrift für Naturforschung*, vol. 41, pp. 1238–1244, 1986.
- [10] N. J. A. Sloane et al., *The On-Line Encyclopedia of Integer Sequences*, [Online]. Available: <https://oeis.org>, Accessed on: November 2021.
- [11] L. I. Stiel and G. Odos, "The normal boiling points and critical constants of saturated aliphatic hydrocarbons," *American Institute of Chemical Engineers*, vol. 8, no. 4, Article ID 5276529, 1962.
- [12] H. Wiener, "Structural determination of paraffin boiling points," *J. Amer. Chem. Soc.*, vol. 69, pp. 17–20, 1947.

Julian D. Allagan

Received May 30, 2022

Accepted June 02, 2022

Julian Allagan
Elizabeth City State University
Elizabeth City, North Carolina, U.S.A
E-mail: adallagan@ecsu.edu

Graph-based decision making for varying complexity multicriteria problems

Oleksandr Nesterenko, Igor Netesin, Valery Polischuk,
Yuri Selin

Abstract

In the modern world in various spheres of activity, the number of problems that need multi-criteria decision-making (MCDM) is constantly increasing. Researchers and experts offer a number of approaches to MCDM process in such tasks; in particular, most of them are based on expert methods. However, in practice, these methods require significant intellectual effort of experts and organizational and technical workload during the expert survey, and also usually take a long time. At the same time, it is not always possible for experts to use certain characteristics of alternatives, which also carries the risk of making decisions based on unfounded expert opinions. Therefore, such methods and tools should be clear and informative and at the same time easy to use to ensure the efficiency and effectiveness of their use.

We offer a graph-based approach to expert decision-making and information visualization processes that meets these requirements and can significantly improve the efficiency of decision-making in multi-criteria selection tasks.

Keywords: information technologies, MCDM, ontologies, expert voting, analytic hierarchy process, analytic network process, graphs, visualization.

1 Introduction

Solving multi-criteria decision-making (MCDM) tasks is largely related to expert intelligence in the field of alternatives. There are many expert methods, but most of them have a number of disadvantages. For

example, the Delphi method is often proposed, but it produces a significant intellectual, organizational and technical burden on the organizers of the survey, causes a large number of iterations in the work of experts and usually requires considerable time to process questionnaires. One of the common methods for ranking alternatives according to certain criteria is the analytic hierarchy process (AHP) mentioned in many works, but it is also not without its drawbacks. Therefore, many experts and researchers offer different approaches to modifying the processes of using common expert methods to overcome existing shortcomings and increase the effectiveness of decisions.

Keep in mind that solving any complex multifactor problems is impossible without modeling and data analysis. Whatever methods are used to evaluate alternatives, in order to support expert decision-making in today's complex information space, it is necessary to ensure the collection, presentation, and analysis at various levels of a significant body of heterogeneous data. At the same time, it is emphasized that the processing of the necessary data is now difficult to imagine without the appropriate means of visualization of information.

Based on this, a set of methods and tools integrated to achieve this goal is needed to properly support expert decision-making. In this approach, we implement a support system that provides a visually interactive interpretation of the three main stages of decision-making – problem analysis, development of alternatives, and their comparison and selection – covering tasks of varying complexity. For each type of task, our system implements visual access to the model, in-depth analysis of generated solutions, and comparison of alternative solutions. Finally, we evaluate the usefulness and ease of use of our system in the field of security.

2 Literature survey and problem statement

In many cases, MCDM support model-driven computational methods. Providing intuitive access to these methods is crucial. Widely used tools to improve understanding of problems and, ultimately, to improve decisions include graphical and information visualization tools. Many studies in the social sciences confirm this conclusion [1]. Studies show

that the data visualized in graphs require less cognitive effort in interpretation, contribute to the effectiveness of communication, clarity, speed, and understanding of complex concepts. Research also examines not only how visualizations convey complex information, but also how to use visualizations in the learning process, for example, in relation to data structures and algorithms [2].

Visualization began with the development of general recommended visualization systems in data analysis processes, which usually illustrated the design of a data set, but could not recommend target results. Therefore, researchers have begun to conduct research in the direction of approaches to visualization, focused on the task of analysis, with modeling of user needs [3, 4]. In fact, these were the first attempts to use visualization capabilities to decision support.

The use of visual representations provides the analyst with an effective method of sifting through a huge amount of information and making informed decisions on critical issues. The paper [5] investigates the impact of information complexity on situational awareness, measured as the density of the graph. The authors claim that the visual signal of the line thickness is an informational value associated with improving time savings and reducing the mental load on the analyst.

However, despite the prevalence of visualization in research and practice, results from different subject areas are rarely shared, although visualizations and their use may be based on general principles. The authors of [6] proposed an integrative model to provide inter-domain support, based on models of understanding visualization and the so-called dual decision-making process. An interactive visualization tool to support multi-criteria decision-making tasks based on the mental model of the user is proposed in [7]. In an environment where the analysis of “human-loop” data is required, covering not only many attributes of alternatives, but also contextual information (domestic policy, customer requirements, cost-effectiveness), the authors’ approach allows users to intuitively explore different criteria and find solutions.

Although data visualization is crucial to help in decision-making, this tool places high demands on the volume, speed, and veracity of data. There is a need for qualified database experts. This is especially important in the case of processing unstructured data from mass

sources, when decision-makers must be able to observe large graphs of visualization [8]. In response to these challenges, the article [9] discusses methods that make data visualization more efficient and effective by directly engaging users who specify their requirements for creating visualizations. In the article [10], the proposed prototype in the mode of comparison of alternatives displays a graph of parallel coordinates, which demonstrates the advantages of experts. To provide high-level summaries of large datasets “at a glance”, heat maps are used, arranged in a grid as tabular histograms with a color mark.

Visualization issues are addressed to support sustainable decision-making in various areas, including administrative management, where in-depth analysis of societal issues and possible policy options is needed. An example of the inclusion of information visualization in the policy analysis process is provided by [11]. Paper [12] proposed a tool to support decision-making based on timeline and taxonomies visualization to manage the capabilities of the defense order portfolio.

In the current trend where information systems are becoming more intelligent, a variety of representations of formal models of context, including graphics, are used in decision-making processes. The aim of the article [13] is to propose a tree-like view of decision-making practices in a contextual graph based on the Contextual Graph formalism. At the same time, ontology-based models occupy a special place among context models. The analysis of the context of the business operation of employment using the context graph was carried out in [14]. For each business operation, its contextual ontology is determined, which reflects the contextual knowledge. Such ontology identifies situation-relevant entities, relationships, and rules. The ontological scheme consists of a hierarchical data structure, contains information about the properties, as well as the relationship between the concepts and objects of the subject area. It is important that the ontology supports decision-making through the possibility of program-interpreted computer representation of knowledge. As a result, it adds intelligence to relevant information technologies in various fields [15].

Most multi-criteria tasks can be represented by hierarchical systems. One of the common expert methods that is well suited for hierarchical data structures and offered in many works is analytic hierarchy

process (AHP) [16]. At the same time, it should be noted that AHP does not lack certain shortcomings, in particular in terms of sensitivity to the clarity of the list of alternatives and limitations. It is also usually necessary to minimize the shortcoming associated with the relationship of consistency as an indicator of the quality of expert assessments. Due to this, the search for the method of multi-criteria analysis that is best suited to solve the problem is often extended either by modified AHP or other methods, as well as the use of ontologies [17, 18].

The main conclusion of the analysis is that such approaches allow finding acceptable solutions only if the state of the subject area is clearly defined, and the experts should be sufficiently qualified specialists. Many studies do not take into account the specifics of evaluating alternatives, due to the fact that expert groups usually include officials who find it difficult to navigate the evaluation methods. At the same time, building models based on the integration of concepts and objects in graphical form still remains a confusing problem in determining the priorities of information support of the decision-making process. All of this suggests that it is advisable to conduct research on further improvement of the typical expert decision-making process in multi-criteria problems of different levels of complexity [19, 20] based on the representation of models in the form of graphs and their visualization.

3 Research on the use of graphs to decision-making

3.1 Decision making and complexity of problems

Decision-making is a complex process that takes at least three consecutive steps: 1) to analyze the problem to be solved, 2) to develop alternative solutions and 3) to choose the best solution. Thus, the problem of decision-making can be formally defined by the scheme $\{X \rightarrow A, \Phi\} \rightarrow a^*$, where X is the set of data representing the problem area, $A = \{a\}$ is the set of alternatives (objects of choice), which can be discrete and continuous; Φ – the principle (function) of choice, according to which, using certain criteria, the advantage in the set of alternatives A is established; and a^* is the chosen alternative (or sev-

eral), which is considered the “best”. There are usually three possible types of decision-making tasks:

- 1) the problem of optimal choice – if the sets X and A are unambiguously defined (fixed), and the principle of choice is formalized;
- 2) the problem of informal choice – if X and A are defined, but Φ cannot be formalized;
- 3) the general problem of decision-making – if X and A do not have defined boundaries (can be supplemented and modified), and Φ is informal.

Tasks of the second and third types are unstructured (poorly defined). Such problems are very difficult (and sometimes impossible) to describe in formal language to give the appearance of the optimal choice problem and they are usually solved by expert methods. In terms of complexity, such tasks can also be classified as simple, complex, and very complex. To reflect the differences between these levels of complexity, it is advisable to use the representation of X and A in the form of oriented graphs (see Fig. 1).

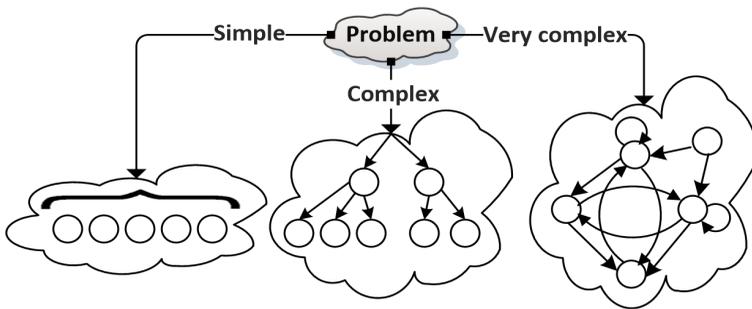


Figure 1. Classification of unstructured problems in terms of complexity

Simple problems can be represented by a linear scheme of alternatives, which in the process of finding a solution are ranked in order of preference over each other. For complex problems that differ in a large number of criteria and characteristics, the search for a solution according to the previous scheme does not give good results.

Usually, such problems are represented by a hierarchical scheme of “criteria – alternative” and require the use of certain algorithms for pairwise comparisons of alternatives according to these criteria. Then, the appropriate calculations are performed on the basis of scalar convolutions of the obtained estimates, taking into account the weights of the criteria.

A hierarchical structure is no longer enough to solve more complex problems. First, for adequate modeling, it is necessary to take into account more parameters of subject areas – objects, factors, requirements, conditions, characteristics, properties, criteria, etc. Second, these parameters can affect each other, and it is important to consider the degree of influence. In this case, it is advisable to use network structures in which the elements of the upper levels may depend on the elements of the lower levels, as well as elements of one level may depend on each other. The network structure allows you to more accurately reflect the relationship in such a subject area. In this case, the elements of the network can be not only simple elements, but also complex elements (components), which in turn consist of a group of homogeneous simple elements. This makes it possible to include in the review almost any knowledge and judgments that may influence the decision.

As the complexity of the problem $\{X \rightarrow A, \Phi\} \rightarrow a^*$, we can take the number of connections (relationships) between the elements of the structure of its model. Denote by $C^{(L)}(X_1, A)$ the complexity of the linear structure, $C^{(I)}(X_2, A)$ – the complexity of the hierarchical structure, and $C^{(N)}(X_3, A)$ – the complexity of the network structure. Then, applying graph theory, these entities can be defined by the following expressions:

$$C^{(L)}(X_1, A) = 0;$$

$$C^{(I)}(X_2, A) = N - 1,$$

where N is the number of vertices in the corresponding oriented tree of criteria; $N - 1$ is the number of edges in this tree;

$$N \leq C^{(N)}(X_3, A) \leq N^2,$$

where N is the number of vertices in the network of components; $N^2 = 2N(N - 1)/2 + N$ – the maximum possible number of arcs and loops in such a network (i.e., in a complete Berge graph – oriented graph without multiple loops and multiple arcs of one direction).

As you can see, the complexity of the network model increases very quickly as new connections are added between its elements. This must be taken into account when building a model for solving a multi-criteria problem.

The procedure of expert formation and evaluation of alternatives is based on the principle of individual and collective work of experts when forming a group of experts, and they can choose from different alternatives using their informal Φ_i . That is, the choice usually depends on the personal preferences of the expert. Thus, overcoming the problem of complexity of tasks also has a negative impact on the subjective vision of experts, which often leads to the preparation of unreasonable decisions.

One of the approaches to solve this problem is a comprehensive information representation of the subject area using a conceptual scheme in the form of ontology, consisting of a hierarchical data structure, containing information about the properties and relationships between concepts and objects of the subject area. As you know, in the general case, computer ontology is formally represented by an ordered trio $O = \langle X, R, F \rangle$, where X is the set of concepts (concepts, terms) of the subject area, R is the set of relations and properties between them, F is the interpretation function (definitions) X and/or R . Finally, as mentioned above, one of the popular tools to improve understanding of the problem and, ultimately, to make effective decisions is to visualize information. Thus, for unstructured tasks, the cognitive decision-making process on an information basis can be presented in Fig. 2.

That is, the choice usually depends on the personal preferences of the expert. At the stage of problem analysis, the collected data is studied and on their basis the initial list of alternatives – “long list” (LL), usually with a linear structure – is determined. At the stage of developing alternatives, it is necessary to select a short list (usually no more than five) – “short list” (SL). It is believed that human thinking is better suited to assessing preferences on multiple objects than on

multiple sets of characteristics. But the advantage of the first approach is only when evaluating fairly simple objects.

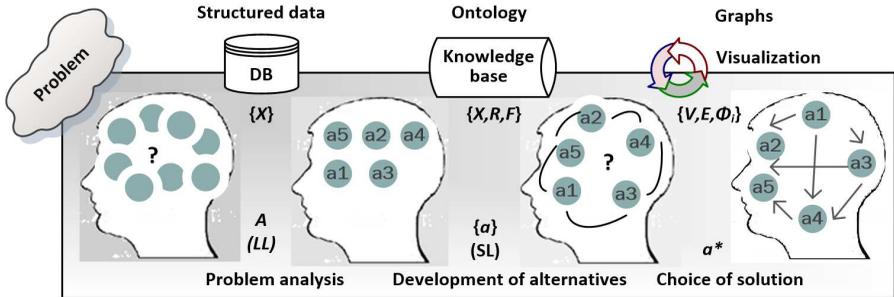


Figure 2. Cognitive decision-making process on information basis

For complex experts, it is easier to determine which of the alternatives is better, given its individual properties (characteristics) – of course, provided they exist. The information space, which should provide experts with comprehensive and clear support of their activities on an objective basis, is formed by the ontological system of the knowledge base. The properties of ontology objects can be used by experts as criteria against which experts can choose alternatives from a variety of possible alternatives.

For complex problems (with a hierarchical structure), the evaluation and selection of the best alternative to SL is based on the principle of direct dominance (greater influence, greater advantage, greater probability), according to which structural elements are compared in pairs (usually on a qualitative scale of linguistic variable). After completing this step, it is necessary to return to the evaluation and comparison of alternatives in general – that is, to perform a composition of criteria. This allows you to find the best of the alternatives or rank them. The ontology of the subject area should clearly define all the characteristics of the criteria, prevent inconsistencies in the selection results due to subjective views, lack of knowledge and errors of experts, and the influence of various factors on them. But such informational support for solving selection problems is not always enough. For example, when

using AHP, the number of spreadsheets, which depends on the number of alternatives, characteristics, and experts, is usually quite significant. This is especially true for recalculations in the event of significant inconsistencies. In addition, the AHP does not check transitive consistency. In order to prevent matrix inconsistencies, it is necessary to “direct” experts in a certain direction in order to avoid extreme subjectivism.

As such a tool, it is proposed to visualize the process of pairwise comparisons in the form of an oriented graph (V, E) with synchronous control of transitivity. The vertices V of the indicated graph correspond to the alternatives, and the edges E with the arrow indicate the advantages of the alternatives. For example, an arc (a_r, a_q) will go from vertex a_r to vertex a_q if $(a_r \succ a_q)$, where the symbol \succ means general superiority. To improve understanding, if necessary, the arcs are loaded with numerical values that correspond to the expert qualitative values of the degree of superiority of one alternative over another according to a certain criterion. Based on their non-formalized Φ_i , the expert can adjust the directions and loads of the edges. Because all selected alternatives are compared in pairs, all vertices will be connected by arcs at the end of the procedure. The resulting graph will be a complete oriented graph, which in graph theory is called a tournament. An evaluation option that satisfies the conditions of transitivity is a must. This stems from the possibility of strict linear ordering vertices of the transitive tournament in the order of their reachability, as all its vertices have different input and output degrees of arcs.

3.2 Complexity of problems and methods of decision making

Practice shows that the most common way of collective decision-making in expert groups is voting. Voting procedures, even if they seem simple, are complex and sophisticated ways of deciding on the basis of conciliation of interests. Finding such a decision is facilitated by the correct choice of the voting procedure, which is characterized by the following stages: a) each participant in the procedure forms his/her opinion on alternatives and reflects it in accordance with the instructions; b) in accordance with one or another formal procedure for processing this

information, a collective decision is determined.

There are numerous voting procedures. Given the above requirement of simplification for experts of the process of forming and evaluating alternatives and choosing the most acceptable one, when these alternatives are fairly simple objects for the basic method, the technique of approval voting (with modification) is proposed. Each expert can both submit his/her proposal for inclusion/exclusion in a variety of alternatives, and participate in the process of improving the proposals of other experts. The main thing is that the expert has the right to support not only one, but also several alternatives, which allow experts to make decisions being closer to consensus than other methods.

But by voting, it is possible to achieve acceptable results only in the case of simple tasks. More complex multi-criteria tasks are usually represented by a hierarchical system. At its lower level, alternatives are evaluated using a vector of criteria formed by the decomposition of properties. At the upper level, with the help of the composition mechanism, the assessment as a whole is formed. One method that is well suited for hierarchical data structures is AHP. In AHP, the hierarchical structure of the problem of choosing alternatives is a graphical representation in the form of an inverted tree. In this structure, each element, except the top, depends on one or more elements above.

For even more complex problems from these classes, when indirect dominance is used to determine the relationship between their elements (alternatives, criteria, characteristics, factors, conditions, scenarios, etc.) network models are needed. To build such models, it is advisable to use the analytic network process (ANP), which is a development of AHP [21]. The ANP involves the construction of an oriented graph without multiple loops and multiple arcs of one direction (Berge graph) and a super-matrix of influences between simple elements and components of the graph. In a super-matrix (block matrix) formed on the basis of a graph, each block is a matrix of pairwise comparisons M_{ij} , which determines the influence of the elements of the i -th component on the elements of the j -th component. After formation of all necessary matrices with application of the corresponding matrix transformations, the algorithm of their calculations is realized to obtain the generalized numerical values. Based on them, the ranking of

alternatives is carried out. Finally, visualizing the process of pairwise comparisons of alternatives in the form of an oriented graph is an additional means of improving the consistency of expert judgments. Thus, the decision-making process, in general, takes place according to the algorithm shown in Fig. 3.

It is not uncommon for a problem, originally defined as simple, to be not quite as it seemed in the simulation process. Therefore, the algorithm provides transitions from one level to another and cyclic return in case of unsatisfactory results.

4 The decision-making process on the example of a typical multi-criteria problem

Let's consider the application of the proposed approach on the example of solving the problem of rating alternatives for extinguishing forest fires (FF), which are the most common dangerous event. This is a typical multi-criteria task facing the organizational unit of the Civil Defense Force whose goal is to determine the composition of means and resources that will have the necessary capabilities to perform tasks, taking into account the importance of tasks and other criteria.

The modern way of extinguishing FF is to involve aviation. The aircraft flight modes during the discharge of fire-extinguishing liquid depend on many factors: the distance from the aerodromes of the permanent base of the aircraft, fire characteristics and level of smoke, the length of the section on the combustion front, and others. Forest fires are accompanied by high combustion temperatures, intense air turbulence and smoke.

Thus, to calculate the forces and means of FF liquidation, the fire extinguishing manager must operate with the values of many data. The formation of the database should take place in advance on the basis of experience and knowledge of experts, taking into account possible situations. An important factor in the presentation of such knowledge is the ontological descriptions of the set of concepts, objects, connections and processes due to the characteristics of this area, which are usually formed using graph models.

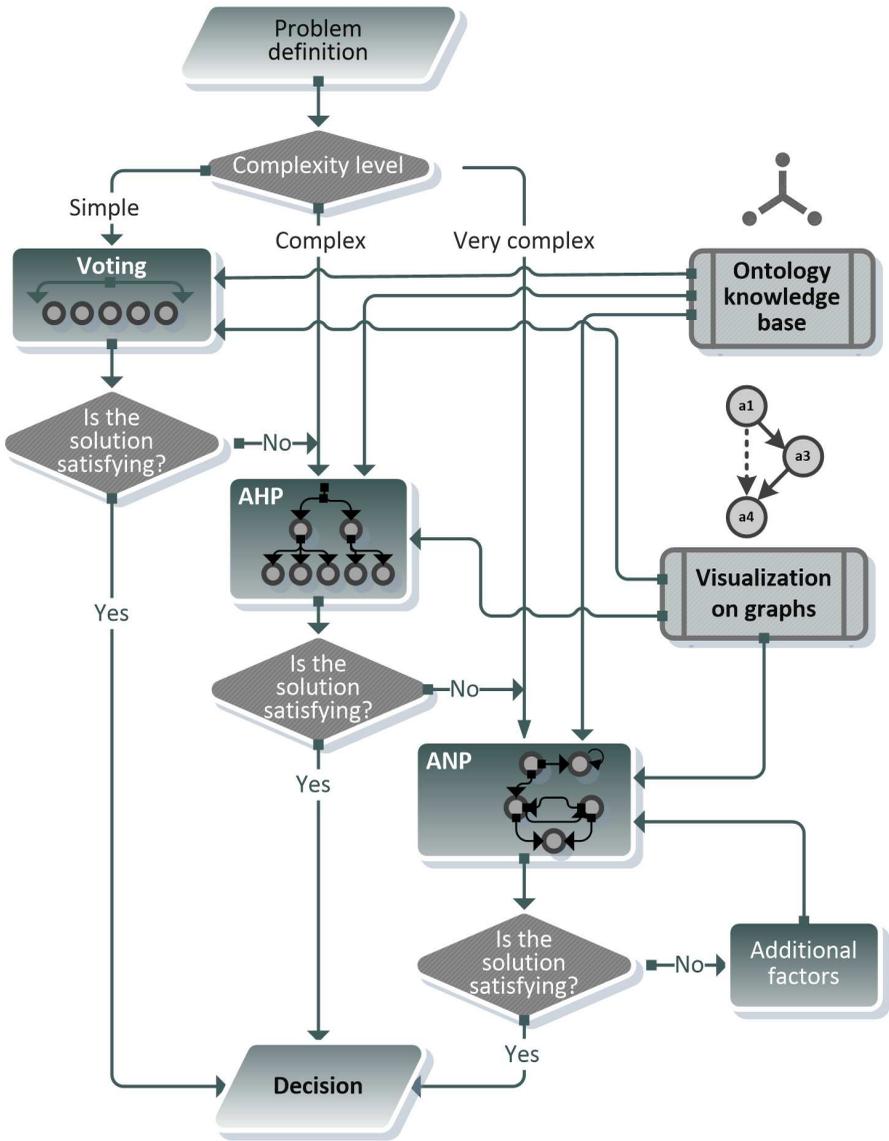


Figure 3. Algorithm of the decision-making process

Clarity of actions on FF liquidation is achieved by development and delivery to extinguishing participants of the aircraft flight schedule. The source information for the schedule compilation should be the values from the database and knowledge base. Based on these data, it is possible to make a significant number of options for aviation tasks (alternatives), among which you need to choose the best for flight schedule preparation.

Thus, according to the developed variants of tasks, five alternatives (C1 – C5) have been proposed for consideration, which can be used for the flight schedule. To move to the evaluation, the necessary characteristics are selected from ontologies of the knowledge base (see Fig. 4).

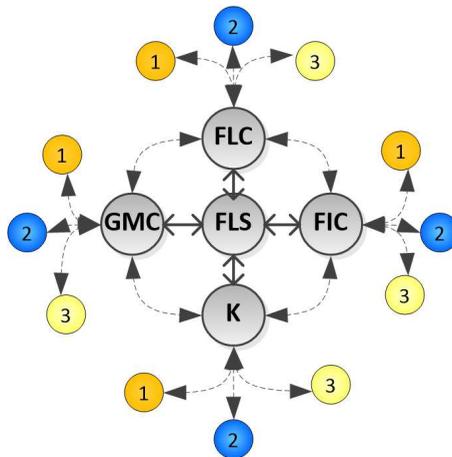


Figure 4. Top level graph of ontologies (FLS – flight schedule; FLC – flight characteristics; FIC – fire characteristics; GMC – characteristics of ground means; K – criteria; 1,2,3, and so on – groups of characteristics)

Initially, all experts vote, using a built ontology graph, which helps them to evaluate the alternatives. An example of display for one of the voting variants for one expert is shown in Fig. 5.

Virtually every voting procedure can lead to the choice of more

Initial rating of alternatives		
Tools 	Add 	Help 
Expert: Igor		Changed: ↓
Isupport	 	23.09.2021
exclude		23.09.2021

Figure 5. Evaluation results by voting procedure

than one alternative or not even being able to identify any alternative. According to the approach under consideration, it is proposed to clarify the decision by transferring the problem to a higher level of complexity and applying a hierarchical method of evaluation.

Using ontological data, a hierarchy of the problem is built. This hierarchy for our model example has three bushes (root subtrees) of criteria, each of which has its own branches. For example, three groups of criteria can be defined: K1 – compliance with the task; K2 – risks in task performing; K3 – cost of flights. The first group includes such criteria as versatility, probability of detection, mobility, reliability, efficiency, range, duration of action, etc. The second group includes the dependence on fire intensity, dependence on wind force, etc. The third group includes the estimated costs of using different aircraft depending on their types, bases, etc. The values of these criteria (characteristics) are presented in databases. The hierarchy for this example is shown in Fig. 6.

According to the AHP algorithm, a unified set of tables is formed for experts to record the results of pairwise comparison of alternatives for each criterion. As a result of processing of tables the standardized values of estimations by all experts of all alternatives in comparison with others on each criterion are determined. After that, they are folded. According to the proposed approach, instead of filling in the tables, the expert compares any pair of alternatives using a special

graphical interface, which displays the vertices of the graph with the names of the alternatives.

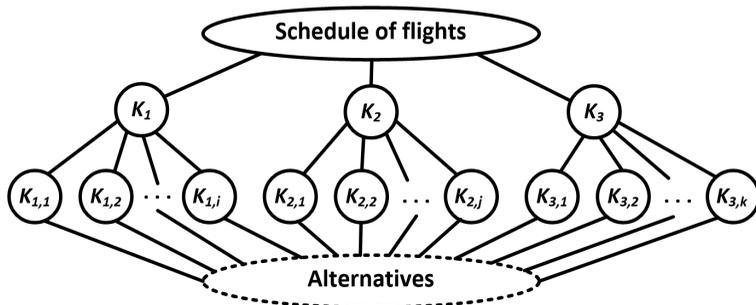


Figure 6. Hierarchy of the task

In this case, tables with the corresponding quantitative estimates (degrees of preference) are formed automatically. If necessary, the expert can review some of his/her own preliminary judgments in order to improve their consistency by editing the graph with preservation of transitivity. An example of the comparison steps ($S1 - S4$) of one of the experts leading to the final transitional tournament is given in Eq. (1), and the visualization of the tournament is shown in Fig. 7.

$$\begin{aligned}
 S1.C2 &\sim C5.\{\{S2.C1 \succ C3.S3.C3 \succ C4.\} \implies \\
 &\implies (C1 \succ C4).S4.C4 \succ C2\} \implies \\
 &\implies (C1 \succ C2, C3 \succ C2).
 \end{aligned}
 \tag{1}$$

In emergencies, alternatives often need to take into account the various elements and entities and the relationships between them. It is often not possible to describe all the necessary relationships in a hierarchical structure. In this case, it is proposed to use a network model based on ANP.

An example of such a more complex task can be the case where the interaction of fire crews with ground rescue units is taken into account in firefighting. To do this, in addition to the general situation, it is

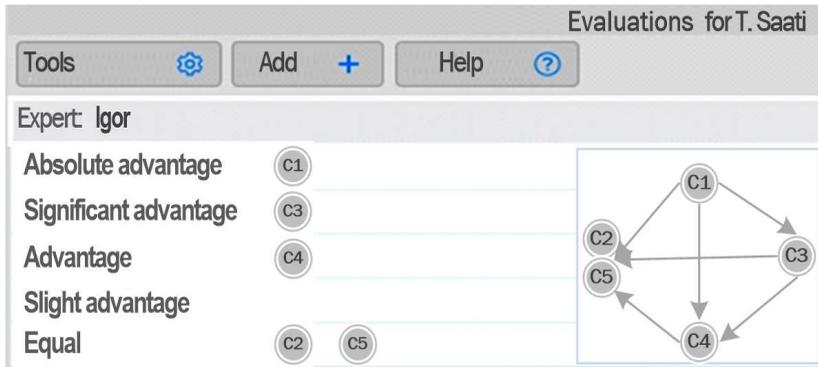


Figure 7. Preliminary visualization of AHP tournament for one of the experts

necessary to take into account various factors and their mutual influences to ensure ground methods of firefighting and security of forces and means. Therefore, when building network models, much attention should be paid to the development of network structure, which should provide the ability to have simple elements (individual entities) and complex elements (components), which in turn consist of simple elements, as their vertices (nodes). Both external dependencies between components and internal dependencies between elements within one component must be taken into account. For example, the ability to maintain the wetness of the local area band by ground units (L1) can significantly reduce the requirement to maintain the wetness only by aircraft, and the involvement of special ground equipment (L2) can affect the number of flights and, consequently, the cost of the operation.

It is advisable to build network models by expanding the already built hierarchical models. With this in mind, the network structure for modeling the interaction of fire crews with ground rescue units is shown in Fig. 8.

An example of the resulting graph of comparisons of one of the experts on one of the criteria is shown in Fig. 9.

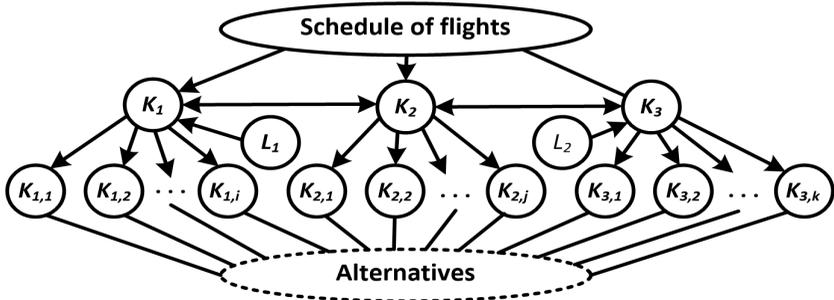


Figure 8. Network structure of the task

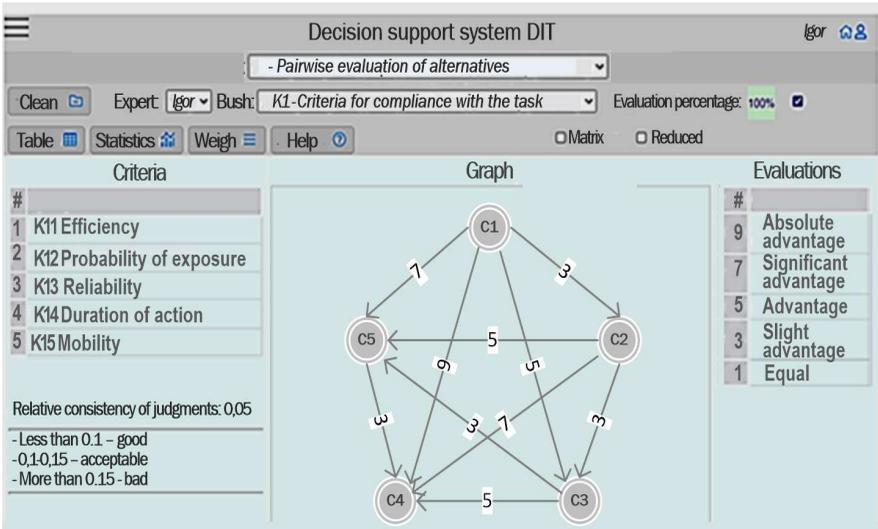


Figure 9. An example of the resulting graph of comparative evaluations for AHP/ANP according to one of criteria of one of experts

5 Conclusions and future work

The proposed approach represents one of the innovative tools for achieving goals and objectives in decision making, which is always relevant. The results of the study are related to the use of graphs as an integrated means of combining into a single set of ontological data models, means of visualizing alternative comparison processes and known methods of multi-criteria analysis. Features of the proposed technique and the results obtained in comparison with existing ones have several advantages. First of all, the proposed approach uses the psychological ability of any person to effectively compare in the presence of visual images. By visualizing on graphs all stages of the process of evaluating and supporting the opinions of experts, expert activity is significantly simplified. In particular, it helps to increase transitive and cardinal coherence. The graphical interface reduces subjectivity and generally creates the conditions for impartiality and fairness. This feature ensures the efficiency, versatility, and simplicity of technical implementation of the decision support procedure.

It should be noted that this study has certain limitations. First of all, they are related to the possibility of building a correct and adequate ontological model of the subject area, which largely depends on the validity and objectivity of the decision. It is necessary to have comprehensive data on the subject area, terminological dictionaries and technical reference books in the electronic presentation, from which it is possible to build an ontological base. At the same time, there is a need to involve qualified specialists in the field of Data Scientist. This is especially true for ANP. In practice, these conditions may not always be met. Usually, when collecting additional data, in particular of a special nature, you may encounter organizational difficulties.

Further directions of this study can be directed on more detailed extension the presented solution of the problem based on the network model. The positive effect of using the potential of this approach may be related to the improvement of the ontological model of the subject area. Given the universality of the approach, the development of this study may consist in its application in various fields.

References

- [1] K. Eberhard, “The effects of visualization on judgment and decision-making: a systematic literature review,” *Manag Rev Q.*, 2021, DOI: 10.1007/s11301-021-00235-8.
- [2] P. Perháč and S. Šimoňáku, “Interactive System for Algorithm and Data Structure Visualizations,” *Computer Science Journal of Moldova*, vol. 30, no. 1(88), 2022, pp. 28–48.
- [3] Sh. Leixian, Sh. Enya, T. Zhiwei *et al.*, “TaskVis: Task-oriented Visualization Recommendation,” in *EuroVis 2021 – Short Papers. The Eurographics Association*, 2021, pp. 91–95. DOI: 10.2312/evs.20211061.
- [4] D. Shi, Y. Shi, X. Xu *et al.*, “Task-Oriented Optimal Sequencing of Visualization Charts,” *arXiv:1908.02502v1*, 2019. DOI: 10.1007/s11301-021-00235-8.
- [5] M. Mittrick, J. Richardson, D.E. Asher, *et al.*, “A Visualization Technique to Improve Situational Awareness,” in *Proc. of the Conference on Next-Generation Analyst VI*, UNSP 106530L, 2018. DOI: 10.1117/12.2304280.
- [6] L.M. Padilla, S.H. Creem-Regehr, M. Hegarty, and J.K. Stefanucci, “Decision making with visualizations: a cognitive framework across disciplines,” *Cogn. Research.*, vol. 3, no. 29, 2018, 25 p. DOI: 10.1186/s41235-018-0120-9.
- [7] Y. Kiml, H. Jeon, and Y-H. Kim, “Visualization support for multi-criteria decision making in software issue propagation,” in *Proc. of the 2021 IEEE 14th Pacific Visualization Symposium (PacificVis 2021)*, pp. 81–85. DOI: 10.1109/PacificVis52677.2021.00018.
- [8] Z. Boulouard, L. Koutti, N. Chouati *et al.*, “Visualizing Large Graphs Out of Unstructured Data for Competitive Intelligence Purposes,” in *Proc. of the SAI Intelligent Systems Conference (IntelliSys 2016)*, pp. 605–626. DOI: 10.1007/978-3-319-56994-9-41.
- [9] X. Qin, Y. Luo, N. Tang, and G. Li, “Making data visualization more efficient and effective: a survey,” *The VLDB Journal*, no. 29, 2020, pp. 93–117. DOI: 10.1007/s00778-019-00588-3.
- [10] E. Hindalong, J. Johnson, G. Carenini, and T. Munzner, “Towards Rigorously Designed Preference Visualizations for Group Decision

- Making,” in *Proc. of the IEEE Pacific Visualization Symposium (PacificVis)*, 2020, pp. 181–190. DOI: 10.1109/PacificVis48177.2020.51111.
- [11] T. Ruppert, J. Dambruch, *et al.*, “Visual Decision Support for Policy Making: Advancing Policy Analysis with Visualization,” in *Policy Practice and Digital Science*, J. Springer, 2015, 10, pp. 321–353. DOI: 10.1007/978-3-319-12784-2-15.
- [12] E.H.S. Lo, N. Tay, and G.J. Bulluss, “Supporting Force Structure Review through graph visualisation and capability view improvements,” in *Proc. of the 21st International congress on modelling and simulation (MODSIM2015)*, pp. 863–869.
- [13] K. García and P. Brézillon, “Model visualization: Combining context-based graph and tree representations,” *Expert systems with applications*, vol. 99, 2018, pp. 103–114. DOI: 10.1016/j.eswa.2018.01.033.
- [14] I.I. Zavushchak, and Ye.V. Burov, “Using contextual graphs to support employment decision-making,” *Lviv Polytechnic National University Institutional Repository*, pp. 129–135, 2018. [Online]. Available: <http://ena.lp.edu.ua>.
- [15] O. Nesterenko and O. Trofymchuki, “Patterns in forming the ontology-based environment of information-analytical activity in administrative management,” *Eastern-European Journal of Enterprise Technologies.*, vol. 101, no. 5/2, 2019, pp. 33–42. DOI: 10.15587/1729-4061.2019.180107.
- [16] T. L. Saaty, *Fundamentals of Decision Making and Priority Theory with the Analytic Hierarchy Process*, Pittsburgh: RWS, 1994.
- [17] F.-x. Xu, X.-h. Liu, W. Chen, *et al.*, “An ontology and AHP based quality evaluation approach for reuse parts of end-of-life construction machinery,” *Mathematical Problems in Engineering*, Article ID 3481030, 2018. DOI: 10.1155/2018/3481030.
- [18] A. Groza, I. Dragoste, I. Sincai, I. Jimborean, and V. Moraru, “An ontology selection and ranking system based on the analytic hierarchy process,” in *Proc. of the 16th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, (SYNASC 2014)*, 2015, pp. 293–300. DOI: 10.1109/SYNASC.2014.47.

- [19] O. Nesterenko, I. Netesin, V. Polischuk, and O. Trofymchuk, “Development of a procedure for expert estimation of capabilities in defense planning under multicriterial conditions,” *Eastern-European Journal of Enterprise Technologies.*, vol. 106, no. 4/2, 2020, pp. 33–43. DOI: 10.15587/1729-4061.2020.208603.
- [20] O. Nesterenko, I. Netesin, V. Polischuk, and Y. Selin, “Multifunctional methodology of expert evaluation alternatives in tasks of different information complexity,” in *Proc. of the 2021 IEEE 3rd International Conference on Advanced Trends in Information Theory (ATIT)*, pp. 226–231. DOI: 10.1109/ATIT54053.2021.9678742.
- [21] T. L. Saaty, *Decision Making with Dependence and Feedback: The Analytic Network Process*, Pittsburgh: RWS, 1996.

Oleksandr Nesterenko, Igor Netesin,
Valery Polischuk, Yuri Selin

Received April 12, 2022
Revised June 11, 2022
Accepted June 14, 2022

Oleksandr Nesterenko
International Europe University
Kyiv, Ukraine ORCID 0000-0001-5329-889X
E-mail: aleksandrnesterenkoua@gmail.com

Igor Netesin
Ukrainian Scientific Center for Development of Information Technologies
Kyiv, Ukraine ORCID 0000-003-1236-287X
E-mail: inetesin@gmail.com

Valery Polischuk
Ukrainian Scientific Center for Development of Information Technologies
Kyiv, Ukraine ORCID 0000-0001-6991-0617
E-mail: valery.polischuk@ukr.net

Yuri Selin
Institute for Applied System Analysis
National Technical University of Ukraine
”Igor Sikorsky Kyiv Polytechnic Institute”
Kyiv, Ukraine ORCID 0000-0002-7562-8586
E-mail: selinyurij1963@gmail.com