

Liu Liu, Bao-sheng Wang, Bo Yu, Qiu-xi Zhong. Automatic malware classification and new malware detection using machine learning. *Frontiers of Information Technology & Electronic Engineering*, 2017, **18**(9):1336-1347.  
<http://dx.doi.org/10.1631/FITEE.1601325>

# Automatic malware classification and new malware detection using machine learning

**Key words:** Malware classification; Machine learning; Malware detection; Feature extraction; Gray scale

Corresponding author: Liu LIU

E-mail: [hotmailliuliu@163.com](mailto:hotmailliuliu@163.com)

 ORCID: <http://orcid.org/0000-0002-6523-1454>

# Motivation

- Traditional malware analysis methods based on artificial experience are difficult to detect various variants of malicious code and dig useful information from massive information.
- Traditional anti-virus systems based on signature fail to classify unknown malware into corresponding families and detect new malware.
- Automated malware analysis based on machine learning is gaining more and more attention from researchers.

# Main idea

- Inspired by the image processing technology, we try to combine the analysis of malicious code with the image processing technology, trying to find more efficient malicious code feature extraction method.
- In order to effectively detect new malicious code, we try to improve the clustering algorithm to make it more suitable for high dimensional space.

# Method

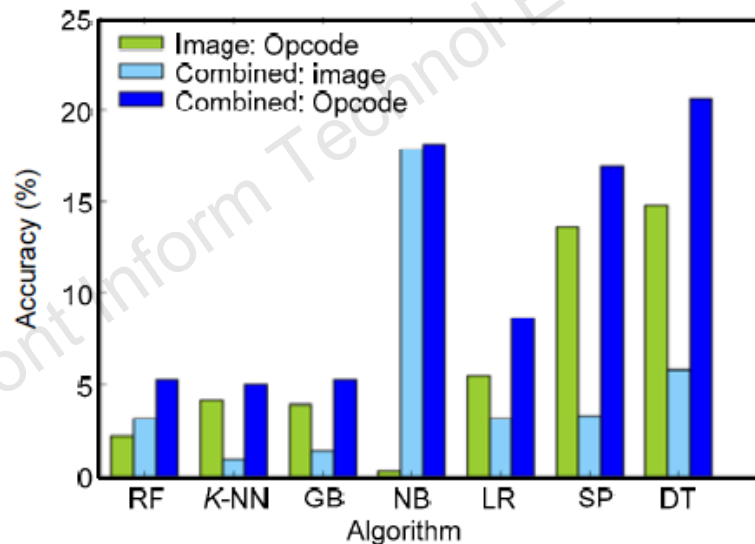
The paper propose an incremental malware detection system, which is able to classify malware families and detect the new malware.

This system is divided into three main parts:

1. It proposes a feature extraction method based on gray scale image, opcode n-gram and import functions.
2. It creates decision-making system to assign unknown malware to corresponding family and screen out suspicious software.
3. It presents an improved shared near neighbors algorithm to mine new malicious code.

# Major results

- Experiment 1 shows that our feature extraction method has a higher accuracy rate on malware classification compared with other methods.

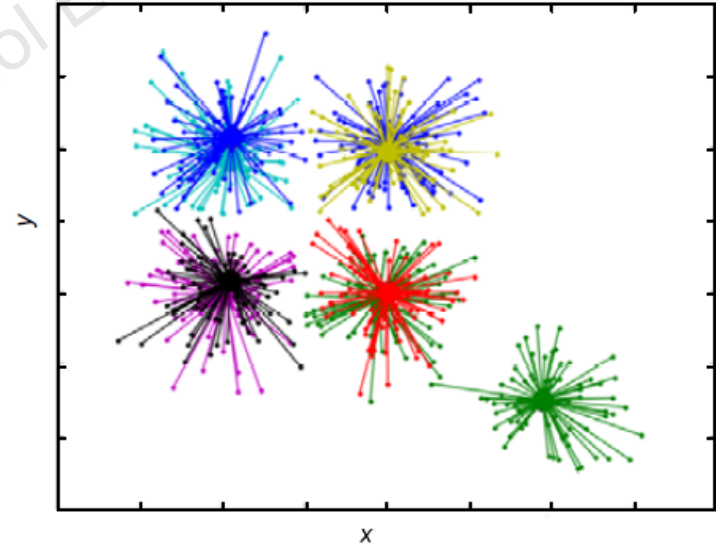
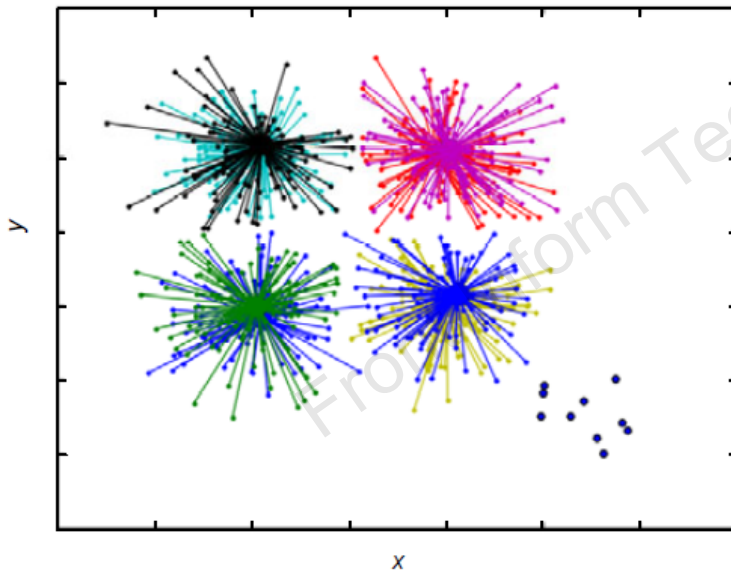


**Fig. 11 Improved accuracies of different classifiers**

RF: random forest; K-NN: K-nearest neighbor; GB: gradient-boosting; NB: Naïve Bayes; LR: logistic regression; SP: support vector machine-poly; DT: decision tree

# Major results

- Experiment 2 shows that our SNN model can detect new malicious samples. And a new family can be clustered in an incremental manner



**New malware detection(left); new malicious family clustering (right)**

# Conclusions

- We propose the texture of malware which effectively describes the features of the difference programs.
- SNN algorithm can accurately detect abnormal samples in high dimensional space, and through the incremental strategy to explore new malicious family.
- Finally, the results of the experiments show that our method effectively improves the accuracy of malware classification, but also can effectively discover new malware.