# Review: Efficient Characterization and Classification of Malware Using Deep Learning

CSE424: Pattern Recognition
Section: 2

Name: Nabil Hossain Chowdhury — Student ID: 21241006

**The authors present a deep learning-based model to classify malware, addressing the limitations of traditional techniques like static analysis, which struggles with packed or obfuscated malware. To improve accuracy and reduce computational costs, they extracted features at three levels; byte-level, basic, and assembly and developed models based on these features to streamline malware classification.**

## 1.1 Motivation

Different malware detection techniques, including signature-based methods, have been proposed. However, signature-based detection is vulnerable to minor code changes, allowing malware to evade detection. On the other hand, Static Analysis struggles with packed malware. To address these limitations, deep learning has become the leading approach for malware detection and classification, due to its ability to reveal the underlying structure of malware data, resulting in higher accuracy and fewer false positives.

## 1.2 Contribution

The primary contribution of this research is the development of a deep learning meta model that is able to dynamically select the most computationally cheap feature set required to classify a sample. This model aims to optimize the trade off between accuracy and computational cost by determining appropriately when to use basic, byte-level or assembly level features.

## 1.3 Methodology

The first objective was to extract static features from malware samples at three levels: basic, byte, and assembly. Based on these levels, three deep learning models were developed and trained on a large dataset of malware samples, primarily targeting financial institutions. The models were designed to identify the simplest feature set that could still accurately classify each malware sample.

## 1.4 Conclusion

The result concluded that the meta-model is highly effective in balancing accuracy and computational efficiency in malware classification. Additionally, by predicting when to use simple or complex feature set, the model reduce the time required for classification, while also keeping a high accuracy over 90%. This approach makes it beneficial for large scale malware classification, where resources can be limited.

## Limitations

**2.1** While the objective is to require low computational resource, often using complex features for classification, such as assembly level requires significant resources, which can be time consuming, especially for larger datasets.

**2.2** The models are trained on a specific dataset with predefined families of malware, there is a need to include additional malware families and to be able to detect a previously unseen malware. So the models may not perform well when a malware significantly deviates from the data on which models were trained.

**2.3** Given the complexity of deep learning models, there is a risk of over fitting, where the models might perform exceptionally well on the training data but fail to generalize to unseen data.

## Synthesis

**3.1** This research and its techniques opens up a lot of possibilities in cybersecurity. It is possible to expand on the current models to create real-time malware detection systems that can operate in environments with constrained resources.

**3.2** Further research can be done by conducting large-scale studies to test the models' scalability to millions of malware samples across diverse environments, including cloud-based systems.