

I-MAD: Interpretable Malware Detector Using Galaxy Transformer

Miles Q. Li^a, Benjamin C. M. Fung^{b,*}, Philippe Charland^c and Steven H. H. Ding^d

^aSchool of Computer Science, McGill University, Montreal, Canada

^bSchool of Information Studies, McGill University, Montreal, Canada

^cMission Critical Cyber Security Section, Defence R&D Canada, Quebec, Canada

^dSchool of Computing, Queen's University, Kingston, Canada

The official version of this article is published in Elsevier Computers & Security.

ARTICLE INFO

Keywords:

Cybersecurity
Malware detection
Deep learning
Transformers
Interpretability

Abstract

Malware currently presents a number of serious threats to computer users. Signature-based malware detection methods are limited in detecting new malware samples that are significantly different from known ones. Therefore, machine learning-based methods have been proposed, but there are two challenges these methods face. The first is to model the full semantics behind the assembly code of malware. The second challenge is to provide interpretable results while keeping excellent detection performance. In this paper, we propose an *Interpretable Malware Detector (I-MAD)* that outperforms state-of-the-art static malware detection models regarding accuracy with excellent interpretability. To improve the detection performance, *I-MAD* incorporates a novel network component called the *Galaxy Transformer network* that can understand assembly code at the basic block, function, and executable levels. It also incorporates our proposed interpretable feed-forward neural network to provide interpretations for its detection results by quantifying the impact of each feature with respect to the prediction. Experiment results show that our model significantly outperforms existing state-of-the-art static malware detection models and presents meaningful interpretations.

1. Introduction

Malware is software written in order to steal credentials of computer users, damage computer systems, encrypt documents for ransom, and other nefarious goals. Recognizing malware samples downloaded by legitimate users in a timely manner is of crucial importance for users' protection. Signature-based malware detection methods are widely used in antivirus products, but they are limited in recognizing significant variants of existing malware and new malware [57, 19]. There is thus a pressing need to create an intelligent malware detection system that has better generability to capture new malware or nontrivial variants of known malware.

Machine learning-based malware analysis methods [57, 37, 17, 6, 48, 20] can automatically learn common patterns of malware from the feature space that have better generalization ability than manually crafted signatures. However, there are two major challenges for machine learning-based malware detection models.

Interpretability is one of the dominant features for classification models in some domains, such as healthcare and cybersecurity. In cybersecurity, the interpretations can help malware analysts justify the classification results and create a knowledge base of malware samples. Hidden Markov model (HMM) [53, 55] and attention-based recurrent neural network (RNN) [15] have been proposed to provide analyzable or interpretable classification results on sequential data. Linear models such as logistic/softmax regression and Naive Bayes produce interpretable results on vectorial data but usu-

ally yield inferior classification performance than non-linear models such as multi-layer feed-forward neural networks [11]. However, the hidden layers between the input and the logistic/softmax layer make multi-layer feed-forward neural networks lose the interpretability of logistic/softmax regression to directly attribute the impact of each feature. It's still a challenge to keep interpretability as well as classification performance for feed-forward neural networks.

As the workload of malware exists mainly in its assembly code, modelling the assembly code could provide important information for malware detection. However, it is challenging to model the whole assembly code of executables because they are very long sequences. An executable of 1 MB could have hundreds of thousands of instructions. No effective training approaches have been proposed to train such long sequences, and the memory consumption cannot be handled with standard hardware for such long sequences.

Deep learning models have achieved significant breakthroughs in understanding natural language when properly trained on large corpora [41, 20, 42]. *Transformer* [52] based models especially achieve state-of-the-art results in natural language understanding and generation [20, 41, 22, 42, 45, 9]. However, their successful applications are mainly on short text, i.e., sentence-level tasks such as paraphrase detection and sentiment analysis [41, 20], or on short-document texts such as reading comprehension and automatic summarization of news articles [22]. For example, the state-of-the-art sequence model *GPT-3* [9] can process sequences of a maximum length of 2,048 tokens. That makes the transference of the success of existing methods to understanding assembly code a challenge. Apart from the fact that assembly code is too long, the differences between natural language and assembly code in the structure composition and basic units stand as another problem to solve.

*Corresponding author.

✉ miles.qi.li@mail.mcgill.ca (M.Q. Li); ben.fung@mcgill.ca (B.C.M. Fung); philippe.charland@drdc-rddc.gc.ca (P. Charland); ding@cs.queensu.ca (S.H.H. Ding)

ORCID(s): 0000-0001-8423-2906 (B.C.M. Fung)

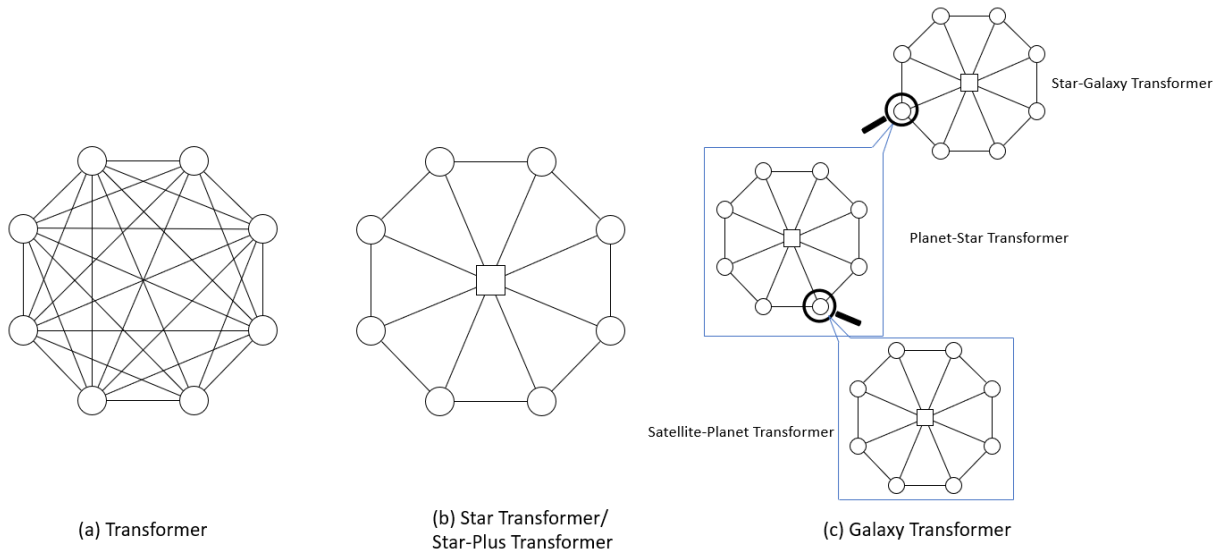


Figure 1: The comparison of topology of the *Transformer*, *Star/Star-Plus Transformer*, and *Galaxy Transformer*.

Despite the fact that the assembly code of an executable is usually very long, it has an innate hierarchical structure: instructions form basic blocks, basic blocks form assembly functions, and assembly functions form the ensemble of assembly code (i.e., the full logic) of an executable. The lengths of basic blocks, assembly functions, and the ensemble of the assembly code of an executable in terms of their direct subunits are usually within thousands. Based on this characteristic, we propose the *Galaxy Transformer* network. It contains three components, namely the *Satellite-Planet Transformer*, the *Planet-Star Transformer*, and the *Star-Galaxy Transformer*. They are three customized *Star-Plus Transformer* networks organized in a hierarchy in order to understand the semantic meaning of the assembly code of an executable at different levels: basic block, assembly function, and executable. The *Star-Plus Transformer* is our improved version of the *Star Transformer* [27], which was proposed for natural language understanding as a variant of the *Transformer* [52]. The time complexity and space complexity of *Transformer* is $O(n^2)$, where n is the length of the token. The *Star Transformer* replaces the fully connected structure of the *Transformer* with a star-shaped topology to reduce the complexities to $O(n)$, and it achieves better results on modestly sized datasets. A comparison of the topology between the *Transformer*, the *Star/Star-Plus Transformer*, and the *Galaxy Transformer* is shown in Figure 1. Our proposed universe-like topology of the *Galaxy Transformer* makes it possible to train very long sequences.

To provide interpretations for the classification results, we propose a novel *interpretable feed-forward neural network (IFFNN)* as the other key component of our full model, the *Interpretable MAlware Detector (I-MAD)*. It has the modelling power of a multi-layer neural network and the interpretability of a logistic regression model. An example of the prediction and its interpretation is given in Table 1. It shows the detection result of a target file, the confidence in the result, the primary contributing features that lead to the prediction, and the most related assembly functions.

The contributions of this paper are summarized below:

1. We propose the *Galaxy Transformer* as an early attempt in the literature to model the full sequences of assembly code for malware detection.
2. We propose two pre-training tasks to train the *Satellite-Planet Transformer* and *Planet-Star Transformer*, which are both components of the *Galaxy Transformer*, to understand the semantic meaning of assembly code at the basic block and assembly function levels.
3. We improve the way to use printable string features and PE import features from previous works with our insights on malware.
4. We propose a novel IFFNN as the classification module of *I-MAD*. It has the same interpretability as logistic regression and the modelling power of multi-layer feed-forward neural networks. It allows *I-MAD* to quantify the impact of each feature for the classification results. It is also a general classification module that can be applied to other classification tasks.

The rest of the paper is organized as follows. Section 2 discusses related work. Section 3 defines the research problem. Section 4 provides the details of our proposed method. Section 5 presents the experiment results and analyses. Section 6 discusses the limitations and our future work. Section 7 concludes the paper.

2. Related Work

2.1. Malware Detection

Malware detection methods fall into three categories: static, dynamic, and hybrid [18]. We summarize the common static and dynamic features in Table 2.

Static methods examine the static content of an executable, while dynamic methods run an executable and analyze its behaviors. Features used in static methods include binary sequences [49, 31, 3, 48, 43, 25], assembly code sequences [37, 17, 2, 3, 47, 25], numerical PE header features [3, 6, 48], PE

Table 1

Sample result of our malware detection and its interpretation, which includes the 5 factors that contribute most to the prediction and the most related assembly functions.

File: 05c199.exe		
Prediction: malicious		
Confidence: 100%		
Primary factors leading to the prediction of malicious		
Feature description	Feature value	Impact
Assembly code	N/A	14.56
Number of PE imports	8	5.12
Major operating system version	1	1.49
Frequency of the string "Sleep"	1	0.82
Frequency of the string ".data"	1	0.59
Most influential assembly functions		
sub_401010		
sub_4010AE		

imports/API calls [49, 6, 48, 39, 25], printable strings [49, 29, 48], and malware images [38, 54, 51]. The most common way to use binary sequences or assembly code sequences is to cut them into n-gram pieces to form features [49, 31, 37, 17, 2, 3, 47, 48]. Some studies find that byte n-grams are effective features [7], while others suggest that byte n-grams are weak or deeply flawed [44] and assembly code is more effective [37, 47]. As control flow graphs could be more robust than assembly code against some obfuscation techniques, there are also instance-based detection methods that identify malware by checking whether an executable contains assembly functions or control flow graphs of known malware [12, 3, 13, 14]. Because they are instance-based methods, they suffer from efficiency issues when the known malware database is large.

Dynamic methods run a target executable in an isolated environment, e.g., a virtual machine or an emulator, and extract features such as the memory image [33, 16, 28], the executed instructions [46, 17, 2, 3], and the invoked system calls or behaviors derived from them [8, 24, 3, 16, 29, 47, 28, 1].

Both static and dynamic methods have their advantages and disadvantages. Compared with static methods, dynamic methods provide more abundant and direct information. Even though both static and dynamic methods extract system calls as features, the parameters passed to those invoked system calls can always be seen with dynamic methods, which is not the case with static methods [8, 16, 29, 47]. Moreover, when a malicious executable is packed or polymorphic, the payload probably cannot be seen by static methods. Yet, to perform its malicious actions it must reveal the payload during execution [8]. This gives another advantage to dynamic over static methods. Therefore, dynamic methods can often achieve better results in the most challenging cases [53]. However, it does not mean that static methods cannot capture malware with those mechanisms, because their use is suspicious and can be detected. Previous works on static malware detection show that when analyzing an unknown executable from multiple feature scopes, it is hard for the malware to

Table 2

Common static and dynamic features for malware detection.

Static	Dynamic
binary sequences	memory image
assembly code	executed instructions
PE header numerical fields	invoked system calls
PE imports/API calls	behaviors
printable strings	
malware images	
control flow graph	

evade detection [3, 29]. On the other hand, one serious shortcoming of dynamic methods is that when malware finds that its execution is being monitored, it may not perform its malicious action to evade detection. Thus, dynamic methods may fail to detect it [8, 57]. In addition, dynamically analyzing an executable is very time consuming.

Hybrid methods extract both static and dynamic features and integrate them into one malware detection model [3, 29, 47, 18]. These two kinds of features are expected to provide complementary information to the model so that it has a more comprehensive view of a sample.

2.2. Transformers

As programming languages and natural languages share some similar characteristics, the experience in modeling the latter can be customized to model the former. Before Vaswani et al. [52] proposed the deep learning model known as the *Transformer*, most state-of-the-art neural machine translation models belonged to the class of *attention-based recurrent neural network (RNN)* models. In these models, an *RNN* is used to encode the source text, and another *RNN* with attention mechanism is used to generate the translation word by word [5, 36]. The attention mechanism is used to determine the importance of the words in the source text for generating each translated word. One disadvantage of this type of model is that the recurrence nature precludes parallelism. Another disadvantage is that the attention mechanism assigns only one importance weight to a word in the source text so it can focus on just one aspect of the words.

The *Transformer* addresses both problems and achieves new state-of-the-art performance on machine translation by abandoning the *RNN* and relying only on an improved attention mechanism [52]. The attention mechanism in the *Transformer* is referred to as *multi-head attention*, which allows multiple attention weights to be assigned to each item. Each weight corresponds to one aspect of an item, thus their attention mechanism is more powerful than the previously proposed attention mechanism in its modeling ability [52]. As there is no *RNN* in it, the recurrence nature of the encoder does not exist anymore, which tremendously increases the parallelism and computing efficiency. Since 2017, researchers have seen the potential of the *Transformer* and proposed their own ways to pre-train the *Transformer* on unlabeled corpora that are abundant and then fine-tune it for downstream NLP tasks. They constantly achieve significantly

better results than previous methods on many NLP tasks [41, 20, 22, 42, 56, 45, 9]. The problem of the *Transformer* is that it computes the attention weights between any two items of a sequence, which leads to $O(n^2)$ time and space complexity. Therefore, the *Star Transformer* [27] is proposed to reduce the complexities to $O(n)$ by adding an additional node to collect global information and connect only adjacent items of the sequence. When the length of a sequence is very long, such a sequential model is still hard to train. For this reason, we propose the *Galaxy Transformer* with a hierarchical topology, so that it has $O(n)$ complexities and can be trained at different levels.

2.3. Interpretable Networks

In most cases, deep learning models are proposed to achieve the best performance for certain research problems without considering their interpretability. However, interpretability is very important in some fields. In healthcare, the rationale for decisions or predictions made by deep learning models and the contributions of different factors leading to them need to be validated by doctors because they concern patients' health [23, 50, 11]. In cybersecurity, deep learning-based malware detection is aimed at replacing signature-based methods to be practical for antivirus products to recognize and then quarantine/delete malware for computer users. However, a deep learning malware detector that cannot explain why an executable is malicious is unlikely to be completely practical. This is because there are false positives, and malware analysts often need to justify detection results. The interpretations of a deep learning-based malware detector alleviate malware analysts' efforts of examining them from scratch and creating a knowledge base of malware samples [11].

Shicke et al. [50] make the criticism that deep learning models are hard to interpret, and therefore linear models dominate applied clinical informatics. They also review some attempts to make deep learning models interpretable. For sequential data, Choi et al. [15] propose the interpretable network *RETAIN* to compute the importance of each variable in patients' medical records to their diagnostic predictions. *RETAIN* is composed of two attention-based *RNNs* to form a softmax regression with dynamically computed weights. For image classification, Zeiler et al. [58] propose Deconvolutional Network (deconvnet) to provide interpretable classification results by revealing which parts of an image are important for its classification. For the classification of vectorial data, logistic/softmax regression and Naive Bayes can interpret how much each feature contributes to a classification result. However, they rely on the feature independence assumption, and thus the interactions of different features cannot be modelled to influence the classification. Inspired by *RETAIN* [15], we propose a novel multi-layer feed-forward neural network to simulate a logistic regression with a dynamically computed weight of each feature to achieve the same interpretability as logistic regression, while keeping the performance of non-linear models.

3. Problem Definition

In this section, we define some important concepts, followed by the definition of the research problem.

An executable is a sequence of bytes:

$$exe = \langle byte_1, byte_2, \dots \rangle \quad (1)$$

The feature set of an executable is extracted by a set of extractors:

$$fea(exe) = \{ext_1(exe), ext_2(exe), \dots\} \quad (2)$$

Except for assembly code, the other extracted features can be represented as a vector. We represent the assembly code as a series of nested sets and sequences.

The assembly code of an executable is a set of assembly functions:

$$code(exe) = \{f_1, f_2, \dots\} \quad (3)$$

An *assembly function* is a set of basic blocks:

$$f = \{b_1, b_2, \dots\} \quad (4)$$

A *basic block* is a sequence of assembly instructions:

$$b = \langle ins_1, ins_2, \dots \rangle \quad (5)$$

An *assembly instruction* is a sequence of one opcode and two operands:

$$ins = \langle Opcode, Operand1, Operand2 \rangle \quad (6)$$

For the uncommon instructions with three operands, the third is ignored. Empty operands are substituted by the special token *EMPTY*. All opcodes and operands form a set, and each of them is assigned an index number. Thus, one instruction can be abstracted as a sequence of three integers, where each integer represents an index of an opcode or operand.

Definition 1 (Malware Detection). *Consider a collection of executables E and a collection of labels L that show the executables in E are benign or malicious. Let exe be an unknown executable that $exe \notin E$. The malware detection problem is to build a classification model M based on E and L such that M can be used to determine whether the executable exe is benign or malicious. ■*

4. Methodology

Our malware detection model *I-MAD* includes the *Galaxy Transformer* to learn a vector to represent the semantic meaning of the assembly code of an executable and an *interpretable feed-forward neural network (IFFNN)* that takes the vector representing the assembly code of a target executable and vectors representing other features as its inputs to produce an interpretable detection result. Figure 2 depicts an overview of our malware detection model. In this section, we introduce the *Star Transformer* and describe how we improve it to form the *Star-Plus Transformer* to build the *Galaxy Transformer*. Then, we propose two methods to pre-train different components of the *Galaxy Transformer*. Next, we introduce the other features we use, our novel *IFFNN*, and how we use it to interpret the detection results.

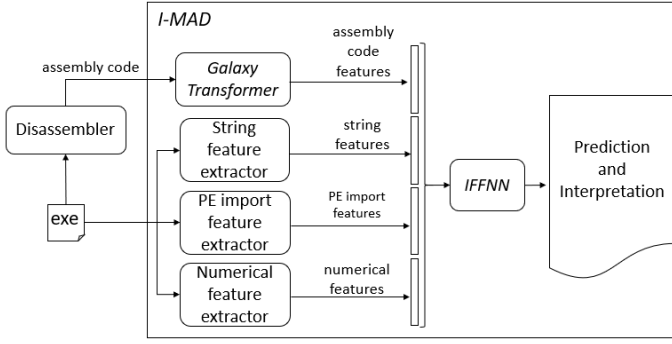


Figure 2: An overview of our I-MAD model.

4.1. Galaxy Transformer

4.1.1. Star Transformer

The *Star Transformer* [27] adopts the multi-head attention from the standard *Transformer*:

$$\begin{aligned} \text{MultiAtt}(q, H) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(qW_i^Q, HW_i^K, HW_i^V) \\ \text{Attention}(q_i, K_i, V_i) &= \text{softmax}\left(\frac{q_i K_i^T}{\sqrt{d_{\text{model}}}}\right)V_i \end{aligned}$$

where $K_i = HW_i^K, V_i = HW_i^V, W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}, W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ are learnable parameters, $q \in \mathbb{R}^{d_k}$ is a query vector, and $H \in \mathbb{R}^{n \times d_k}$ is a matrix that contains vector representations of n items to attend to. To compute the self-attention of a sequence $X = \langle x_1, x_2, \dots, x_n \rangle$, in the *Transformer*, each x_i is a query q , and it attends to all items in the sequence, so $H = X$. Thus, its computational complexity is $O(n^2)$.

To reduce the computational complexity, the *Star Transformer* only considers connections between adjacent items and between a relay node and each item, as shown in Figure 1b. First, for each item x_i , a vector e_i is computed as the summation of its non-contextual semantic embedding and its positional encoding in the same way as the *Transformer* does:

$$\begin{aligned} e_i &= \text{Emb}(x_i) = \text{SE}(x_i) + \text{PE}(i) \\ E &= [e_1; \dots; e_n] \end{aligned}$$

Then, the embeddings are fed into a multi-layer neural network to compute the hidden state for each x_i . h_i^t represents the hidden state of x_i at layer t . h_i^0 is initialized as e_i . The initial hidden state of the additional relay node is $s^0 = \frac{1}{n} \sum_n e_i$. To compute h_i^t , its context matrix C_i^t is formed by the hidden states of itself h_i^{t-1} and its adjacent nodes of the previous layer $h_{i-1}^{t-1}; h_{i+1}^{t-1}$, its embedding e_i , and the hidden state of the relay node s^{t-1} :

$$C_i^t = [h_{i-1}^{t-1}; h_i^{t-1}; h_{i+1}^{t-1}; e_i; s^{t-1}]$$

So, we have $C_i^t \in \mathbb{R}^{6d_{\text{model}}}$.

At each layer, we have

$$\begin{aligned} h_i^t &= \text{LayerNorm}(\text{ReLU}(\text{MultiAtt}(h_i^{t-1}, C_i^t))) \\ H^t &= [h_1^t; \dots; h_n^t] \\ s^t &= \text{LayerNorm}(\text{ReLU}(\text{MultiAtt}(s^{t-1}, H^t))) \end{aligned}$$

Thus, the relay node s^t serves as a global information collector. h_i^t collects local information from its adjacency nodes and global information from s^{t-1} . The computational complexity to compute all h_i^t is $O(n)$, and to compute s^t it is also $O(n)$. The overall computational complexity is therefore $O(n)$.

To put it all together, we represent a *Star Transformer* Layer as follows:

$$H^{t+1}, s^{t+1} = \text{STL}^t(H^t, s^t, E)$$

The full computation of the *Star Transformer* is as follows:

$$\begin{aligned} E &= [\text{Emb}(x_1); \dots; \text{Emb}(x_n)] \\ H^0 &= E, s^0 = \frac{1}{n} \sum_n e_i \\ H^T, s^T &= \text{STL}^T(\text{STL}^{T-1}(\dots \text{STL}^1(H^0, s^0, E), E), E) \end{aligned}$$

4.1.2. Star-Plus Transformer

As previously shown, the *Star Transformer* can generate a contextual vector representation for each item in a sequence and a vector representation for the whole sequence with $O(n)$ computational complexity. We propose the following modifications for better performance.

1. There is no obvious reason why e_i should be in the context matrix C_i^t , so we remove e_i from C_i^t , resulting in $C_i^t = [h_{i-1}^{t-1}; h_i^{t-1}; h_{i+1}^{t-1}; s^{t-1}]$.
2. There was a pointwise feedforward neural network ($\text{FFN} = \text{max}(0; xW_1 + b_1)W_2 + b_2$) after the multi-head attention computation in the *Transformer*, but it is removed in the *Star Transformer* without an explanation for the rationale. We add it back to compose the information collected by all attention heads and to generate higher-level features for the next layer.
3. A max-pooling on H^T across the top layer mixed with s^T was used as the representation for the whole sequence in the *Star Transformer*. We use only s^T to represent the whole sequence, since it has collected global information of the sequence.

To put it together, we have a *Star-Plus Transformer* layer $H^{t+1}, s^{t+1} = \text{SPTL}^t(H^t, s^t)$ computed as follows:

$$\begin{aligned} h_i^{t'} &= \text{LayerNorm}(\text{ReLU}(\text{MultiAtt}(h_i^{t-1}, C_i^t))) \\ h_i^t &= \text{LayerNorm}(\text{ReLU}(\text{FFN}(h_i^{t'}))) \\ H^t &= [h_1^t; \dots; h_n^t] \\ s^{t'} &= \text{LayerNorm}(\text{ReLU}(\text{MultiAtt}(s^{t-1}, H^t))) \\ s^t &= \text{LayerNorm}(\text{ReLU}(\text{FFN}(s^{t'}))) \end{aligned}$$

4.2. Satellite-Planet Transformer to Understand Basic Blocks

As we have stated before, a *basic block* is a sequence of assembly instructions: $b = \langle ins_1, ins_2, \dots \rangle$. The objective of the *Satellite-Planet Transformer* is to learn a vector representation for b using its instructions. To build the *Satellite-Planet Transformer* with the *Star-Plus Transformer*, we modify the input layer of the latter because each instruction is not an atomic item, but a sequence of three items (i.e., an opcode and two operands). Since both the embedding of an instruction ins_i and its positional encoding should have d_{model} dimensions, we make the embeddings of the opcode and operands $d_{model}/3$ dimensions and use the concatenation of them as the embedding of the instruction. It is then added with the positional encoding to form e_i . The concatenation of the vector representation of the opcode and operands to form the vector of an instruction was also previously adopted by Ding et al. [21]. For the output, we directly use s^T , which is the representation of the relay node at the top layer as the semantic meaning representation of the basic block. To train the *Satellite-Planet Transformer* we propose the *Masked Assembly Model* task.

Definition 2 (Masked Assembly Model). *Let (b, ins) be a basic block and assembly instruction pair. Consider a set of basic block and assembly instruction pairs B . For each pair $(b, ins) \in B$, there is one mask instruction m in b that should originally be ins . Let t be a target basic block that is not in any pair of B , and one of its instructions is replaced by m . The Masked Assembly Model task is to build a classification model M based on B to predict the original instruction t replaced by m . ■*

This task is inspired by the *Masked Language Model* task proposed by Devlin et al. [20]. In that task, the authors mask random words from sentences and use the *Transformer* to predict the masked words based on the contextual words in the sentences. Their method is to feed the output vector of the *Transformer* corresponding to a masked word to an output softmax over the vocabulary. The prediction requires both global context and local context. The global context means the semantic meaning of the whole sentence except the masked word. The local context means the position of the masked word and its surrounding words that could indicate what ingredient the missing word should be. As the output vector corresponding to the masked word is the only information source for the output layer to make the prediction, it has to capture both global and local context. This does not fit our objective, since the output vector should only contain the semantic meaning of a basic block (i.e., global contextual information). Therefore, we separate the two kinds of information in two vectors: s^T containing the global contextual information and the output vector of the masked instruction $m = [MASK_OPC, EMPTY, EMPTY]$ containing the local contextual information. We concatenate these two vectors to form one vector and feed it to three feed-forward neural networks with softmax over the whole set of opcodes and operands to predict the opcode and two operands of the

original masked instruction. It should be noted that after this training step, we only need to keep the *Satellite-Planet Transformer*, which generates s^T , the semantic representation of the entire basic block, because the three feed-forward neural networks to predict the original masked instruction are not needed after the training for the *Masked Assembly Model* task.

4.3. Planet-Star Transformer to Understand Assembly Function

The *Planet-Star Transformer* is another customized *Star-Plus Transformer* built on top of the *Satellite-Planet Transformer* to learn the vector representation of the semantic meaning of an assembly function f from the set of vectors representing its basic block $\{b_1, b_2, \dots\}$. As the input is already vectors rather than integers, we abandon the input embedding layer of the *Star-Plus Transformer* that maps integers to embeddings. We directly feed the vectors representing the basic blocks in positional order to form a sequence to the *Planet-Star Transformer*, which is a *Star-Plus Transformer* without an input layer. We use s^T as the vector representation of the assembly function. To train the *Planet-Star Transformer*, we propose the *Assembly Function Clone Detection* task.

Definition 3 (Assembly Function Clone Detection). *Let (f_1, f_2) be an assembly function pair. Let (f_1, f_2, l) be a labeled assembly function pair in which the label l indicates whether the two assembly functions f_1 and f_2 are clones (i.e., semantically equivalent) of each other. Consider a collection of labeled assembly function pairs F . Let $p = (f_1, f_2)$ be a new function pair that p is not any function pair in F . The assembly function clone detection task is to build a classification model M based on F to determine whether the two functions in p are clones of each other. ■*

The intuition is that if the vector representations of assembly functions can be used to determine whether two functions are clones of each other, then they contain the semantic meaning of the assembly functions. We train the network to generate similar vectors in cosine measure (i.e., $\cos(s_{f_1}^T, s_{f_2}^T)$) for real assembly function clone pairs and dissimilar vectors for non-clone pairs. The way we form the function pair dataset is described in Section 5.

4.4. Star-Galaxy Transformer to Understand Full Logic of Executable

Next, we use the *Star-Galaxy Transformer* to learn one vector representing the full logic of an executable based on the representations of all its assembly functions: $\{f_1, f_2, \dots\}$. Technically, this is similar to learning the representation of an assembly function from the representations of its basic blocks, since both are intended to learn one vector representation from a set of vectors. Therefore, the *Star-Galaxy Transformer* is a duplicate of the *Planet-Star Transformer*. Their difference is that they work at different levels of the hierarchy. The representation of the assembly code of an executable generated by the *Star-Galaxy Transformer* is fed to

the *IFFNN* for malware detection without other pre-training tasks proposed for it.

With this, we have completely described how we build the *Galaxy Transformer* with three customized *Star-Plus Transformers* in a hierarchy to compute the vector representation of the assembly code of an executable.

4.5. Other Features

When malware is packed, or is polymorphic or metamorphic, the assembly code of its payload is encrypted and not statically accessible. Hence, using only assembly code would fail to identify its malicious purpose. According to the experience of previous works [3, 29], static analysis can still be effective, because the use of the stealthy mechanisms can be captured when analyzed from multiple static feature scopes. Next, we describe the three kinds of features we use and how we improve the way to use them.

4.5.1. Printable Strings

According to the literature [49, 29, 16, 28], printable strings are important features, because they include, among others, runtime-linked libraries, functions, and registry keys that are commonly used by malware, system paths, and sometimes the names of user-defined functions. Hence, we extract printable strings from the whole byte sequence of an executable. In our algorithm, a continuous subsequence is a *printable string* if it satisfies three conditions: 1) all of its bytes are ASCII characters, 2) it is terminated with a null symbol, and 3) its length is at least 5 bytes. We count the number of instances of each printable string in the training set and put the strings that appear more than a certain threshold, 1,000 in our case, in the frequent string set. Their frequencies in an executable are used as features. This is not new compared to previous works. The improvement is that we also use the number of printable strings that are not in the frequent string set, i.e., uncommon strings, as a feature, and we use the total number of common printable strings in the executable as another feature. This is based on the intuition that encrypted malware has more uncommon printable strings and benign software has more common strings.

4.5.2. PE Imports

PE Imports are dynamically linked libraries and functions shown in the import address table of PE headers. The imports of an executable often illustrate its behaviors, e.g., modify the registry or hook a procedure [49, 48]. The total number of imports show whether the executable is hiding its potential behaviors, because abnormally few imports indicate that runtime linking is largely used or most of its imports are hidden in encrypted data. Therefore, we compute these features in the same way as we compute the printable string features.

4.5.3. PE Header Numerical Features

There are many numerical fields in PE headers that contain information that could form different patterns among malware and benign software (benignware) [6, 48]. Hence,

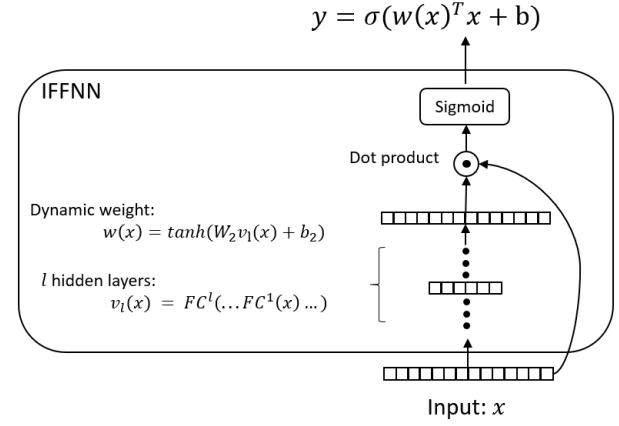


Figure 3: Our proposed *IFFNN*.

we also use these values as features. The list of PE header numerical features we use is given in the supplemental file.

We concatenate the vector representing the full logic of an executable v_{code} , printable string feature vector v_{str} , PE header numerical feature vector v_{num} , and PE import feature vector v_{imp} to form a vector representation of the executable from multiple scopes $v = [v_{code}, v_{str}, v_{num}, v_{imp}]$.

4.6. Interpretable Feed-Forward Neural Network

Interpretability is an important quality of a machine learning model for malware detection. Inspired by the work of Choi et al. [15], which uses two attention-based *RNN* networks to form a softmax regression model with dynamic weights, we propose a novel interpretable feed-forward neural network (*IFFNN*) to form a "dynamic logistic regression" model. Figure 3 illustrates its architecture.

Let $x \in \mathbb{R}^m$ be a feature vector representing a sample. We first feed it to l fully-connected hidden layers:

$$v_l(x) = FC^l(\dots FC^1(x) \dots) \quad (7)$$

$$\text{where } FC^i(v_{i-1}(x)) = \tanh(W_1^i v_{i-1}(x) + b_1^i) \quad (8)$$

where $W_1^i \in \mathbb{R}^{d_h^i \times d_h^{i-1}}$, $b_1^i \in \mathbb{R}^{d_h^i}$, and $v_l(x) \in \mathbb{R}^{d_h^l}$. Then, we apply another normal fully-connected layer of which the output vector has the same dimension as x :

$$w(x) = W_2 v_l(x) + b_2 \quad (9)$$

where $W_2 \in \mathbb{R}^{m \times d_h^l}$, $b_2 \in \mathbb{R}^m$, and $w(x) \in \mathbb{R}^m$. $w(x)$ serves as a weight vector for each feature in x . The final confidence that the input sample is positive (in malware detection, positive means malicious) is calculated as follows:

$$y = IFFNN(x) = \sigma(w(x)^T x + b) \quad (10)$$

$$\text{where } \sigma(z) = \frac{1}{1 + e^{-z}} \quad (11)$$

where $b \in \mathbb{R}$ is a bias term. This is similar to a logistic regression (i.e., $y = \sigma(w^T x + b)$, where w is a parameter vector), and the difference is that our weight vector $w(x)$ is

dynamically computed based on x rather than static parameters.

As can be seen, the *IFFNN* has the same interpretability (a.k.a., intelligibility [35]) as logistic regression because of its decomposability and algorithmic transparency [34] when the features have certain meanings. Since the weight of each feature is dynamically computed by a multi-layer neural network, the feature interactions are still modelled and the weight of each feature to the prediction is contextualized. Thus, it has the modelling power of a non-linear model.

The *IFFNN* can be used for any binomial classification task and can be plainly generalized to a "dynamic softmax regression" model for multinomial classification, as long as the sample can be represented as a vector with fixed dimension.

We feed v , the feature vector from multiple scopes of an executable, to the proposed *IFFNN* to get the confidence y that it is malicious: $y = IFFNN(v)$ and interpret the result. Thus, we complete the full network of the top-level model.

4.7. Attribution

For logistic regression: $y = \sigma(w^T x + b) = \sigma(w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b)$, where $x = (x_1, x_2, \dots, x_m)$ and $w = (w_1, w_2, \dots, w_m)$, the attribution is simple. Whether feature x_j makes the sample positive depends on the sign of $w_j x_j$. If $w_j x_j > 0$, x_j makes it positive, and vice versa. The degree of the impact of x_j on y depends on $|w_j x_j|$: a large $|w_j x_j|$ implies a large impact of x_j . If the model predicts a sample to be positive, the most influential factor that leads to the result is $\max_j w_j x_j$. If the model predicts a sample to be negative, the most influential factor that leads to the result is $\min_j w_j x_j$.

The same idea is applicable to our *IFFNN* for malware detection. Its top layer is logistic regression with dynamically computed weight: $y = \sigma(w(v)^T v) = \sigma(w(v)_1 v_1 + w(v)_2 v_2 + \dots + w(v)_m v_m)$. If $|w(v)_j v_j|$ is large and $w(v)_j v_j > 0$, feature v_j has a large impact on the prediction of malicious. If $|w(v)_j v_j|$ is large and $w(v)_j v_j < 0$, feature v_j has a large impact on the prediction of benign. For printable string features, PE imports, and PE header numerical features, each dimension of their vector representations corresponds to a specific feature. The features can be the frequency of a certain string, whether a certain DLL is imported, the value of a certain numerical field, etc. By checking its $w(v)_j v_j$, we know whether it makes the executable more likely malicious or benign. For the vector representing the full logic of an executable: v_{code} , each of its dimensions has no specific meaning, but we can see the impact of the full logic of the executable by computing the summation of the impact of each dimension of its vector: $\sum_{j \in v_{code}} w_{code,j} v_{code,j}$. If it is positive, from the perspective of the assembly code, the executable is more likely malicious, and vice versa.

As v_{code} is computed by our *Star-Galaxy Transformer* network, the attention weights of the assembly functions to the relay node at the top layer indicate the importance of each assembly function. We compute the summed attention weights of each assembly function over all heads to the relay

node to determine which assembly functions are the main factors that influence the classification results.

4.8. Model Training

To train the *Satellite-Planet Transformer*, the objective function is the cross entropy loss of the prediction on the masked opcode and operands against the real opcode and operands. To train the *Planet-Star Transformer* and simultaneously fine-tune the *Satellite-Planet Transformer*, the objective function is the mean squared error between the computed cosine similarity between two assembly functions and the gold standard (i.e., 1 for clone function pairs, and -1 for non-clone function pairs). To train the full top-level network including the *IFFNN* and the *Star-Galaxy Transformer*, the objective function is the cross entropy loss of the prediction against the real label. To ensure that the *Star-Galaxy Transformer* gets sufficient training, we first train it without concatenating any other feature, i.e., feed v_{code} instead of v to the *IFFNN* ($y = IFFNN(v_{code})$), and train it for malware detection. This is in fact the pre-training of the *Star-Galaxy Transformer*. Then, we concatenate v_{code} with other features to feed it to the *IFFNN* ($y = IFFNN(v)$), and train it the same way for malware detection. The *Satellite-Planet Transformer* and *Planet-Star Transformer* networks are not fine-tuned when we train the top-level network. For all the training objectives, we use Adam [30] with the initial learning rate $1e - 4$. We use early stopping with the validation set to avoid overfitting [10].

5. Experiments

The objectives of our experiments are to 1) evaluate the performance of *I-MAD* for malware detection, 2) compare *I-MAD* to other state-of-the-art static malware detection solutions, and 3) demonstrate the interpretability of *I-MAD*.

We train and evaluate the models on a server with two Xeon E5-2697 CPUs, 384 GB of memory, and four Nvidia Titan XP graphics cards. We use PyTorch [40] to implement our model. We use the "pefile"¹ library to extract numerical features from PE headers. The list of hyper-parameters we tune is given in the supplemental file.

5.1. Datasets and Pre-training

For the two pre-training tasks, we compile several open source projects that are compatible with GCC and/or LLVM. We choose these two compilers because they are the most appropriate options to provide different compilation options to generate semantically equal but literally different assembly functions. GCC compiler provides four different optimization levels (i.e., O0, O1, O2, and O3) to compile projects. We compile *busybox*, *coreutils*, *libcurl*, *libgmp*, *libtomcrypt*, *libz*, *magick*, *openssl*, *puttygen*, and *sqlite3* with GCC at all four optimization levels. Thus, for every assembly function in those projects we have four semantically equivalent versions. O-LLVM² is an obfuscator of the LLVM compiler

¹<https://github.com/erocarrera/pefile>

²<https://github.com/obfuscator-llvm/obfuscator/wiki>

Table 3

Top 10 majority malware families of the dataset.

Malware Family	Number	Percentage
Fareit	9,436	8.2%
Zbot	6,433	5.6%
Emotet	6,343	5.5%
Gandcrab	4,120	3.6%
Mepaow	4,055	3.5%
CobaltStrike	3,151	2.7%
Allaple	2,081	1.8%
Ursnif	1,552	1.3%
Autoit	1,017	1.0%
NaKocTb	794	0.7%
Total	38,982	33.9%

that provides control flow flattening, instruction substitution, and bogus control flow obfuscation mechanisms. We use O-LLVM to compile *libcrypto*, *libgmp*, *libMagickCore*, and *libtomcrypt* with five different settings: no obfuscation, each of the three obfuscation mechanisms, and all three mechanisms. Thus, we have five versions of their every function. We use IDA Pro³, a commercial disassembler, to disassemble our compiled executables and acquire the assembly functions.

We use basic blocks of lengths between 5 to 250 instructions to form our *Masked Assembly Model* dataset; these blocks are within the typical length range of blocks that provide enough context and are not too long to harm training efficiency. As a result, this dataset contains 38,427,440 basic blocks. We use all of them for training and none for testing as the purpose of the dataset is to train the *Satellite-Planet Transformer* to understand assembly code, and the accuracy of this task is uninformative.

We use the semantically equivalent but literally different functions we compiled to form real function clone pairs. We randomly pair the same number of functions to be non-clone function pairs to create the dataset for the *Assembly Function Clone Detection* task. We limit the maximum number of instructions per basic block to be 50 and the maximum number of blocks per function to be 50 in this dataset, so that the memory of our graphics cards can hold the data flowing in the bottom two-level networks. There are 213,656 function pairs in the training set, 26,898 functions in the validation set, and 26,746 functions in the test set. Our bottom two-level networks get a classification accuracy of 91.5% on the test set. This means that the assembly function representations it computes and the representations of basic blocks that are fed to it indeed capture the semantic meanings of assembly code. We do not elaborate on the experiments for this task since it is not the objective task, but rather a task to pre-train the *Planet-Star Transformer* and fine-tune the *Satellite-Planet Transformer*. For malware detection we collected a dataset containing 115,000 benign and 115,000 malicious executables. There is no redundancy in the dataset. Following the literature [49, 31, 43], the benign executables are the

³<https://www.hex-rays.com/products/ida/>

Table 4

Top 10 packers used in the malware dataset.

Packer	Number	Percentage
UPX	7,776	6.7%
BobSoft Mini Delphi	5,262	4.5%
ASProtect	1,826	1.59%
ASPack	1,780	1.55%
PECompact	586	0.51%
Armadillo	369	0.32%
D1S1G	155	0.14%
WinrarSFX	124	0.11%
MoleBox	69	0.06%
WinZipSFX	38	0.03%
Total	17,985	15.6%

.exe and .dll files from the installation paths of software programs. The malicious executables are collected from *MalShare* and *VirusShare*. The top 10 major malware families of the dataset are presented in Table 3. They are obtained with ClamAV⁴. The top 10 known packers that are applied on the malware samples are shown in Table 4. The usage of packers is acquired with Yara Rules⁵. The way we split the dataset into training set, validation set, and test set is introduced in Subsection 5.3.

5.2. Models for Comparison

We compare our *I-MAD* model to several state-of-the-art static malware detection models.

- **Mosk2008OB** Moskovitch et al. [37] propose to use TF or TF-IDF of opcode bi-grams as features and use document frequency (DF), information gain ratio, or Fisher score as the criteria for feature selection. They apply Artificial Neural Networks, Decision Trees, Naïve Bayes, Boosted Decision Trees, and Boosted Naïve Bayes as their malware detection models.
- **Bald2013Meta** Baldangombo et al. [6] propose to extract multiple raw features from PE headers and use information gain and calling frequencies for feature selection and PCA for dimension reduction. They apply SVM, J48, and Naïve Bayes as their malware detection models.
- **Saxe2015Deep** Saxe et al. [48] propose a sophisticated deep learning model that works on four different features: byte/entropy histogram features, PE import features, string 2D histogram features, and PE metadata numerical features. We tried to follow the exact features they extract when we implement it, but they do not provide the exact metadata numerical fields they use, so we just use the same numerical fields of PE headers used in our model as part of their input.
- **Raff2017MalC** Raff et al. [43] treat an executable as a sequence of bytes and apply a gated 1D convolutional neural network (CNN) to classify an executable.

⁴<https://www.clamav.net/>

⁵<https://github.com/Yara-Rules/rules>

Table 5

Results of k-fold cross-validation experiment. It includes the p-values (pv) of t-test for F1 and accuracy between *I-MAD* (ST+) and other models.

Model	P	R	F1	pv (F1)	Acc	pv (Acc)
<i>Mosk2008OB</i>	96.1	95.8	95.9	3.3e-13	95.9	1.6e-20
<i>Bald2013Meta</i>	96.5	95.9	96.2	1.1e-13	96.2	6.7e-20
<i>Saxe2015Deep</i>	95.2	96.1	95.7	4.0e-14	95.6	4.5e-21
<i>Raff2017MaIC</i>	95.9	96.3	96.1	5.6e-15	96.1	4.0e-20
<i>Krcal2018Conv</i>	93.2	93.2	93.2	1.7e-15	93.2	1.0e-23
<i>Mour2019CNN</i>	72.6	71.5	72.0	2.3e-26	71.8	1.5e-30
<i>SVM</i> (same features)	96.1	96.4	96.2	3.7e-13	96.2	5.6e-20
<i>I-MAD</i> (no code)	96.5	96.6	96.5	9.8e-13	96.5	4.0e-19
<i>I-MAD</i> (ST)	97.0	97.9	97.3	5.0e-3	97.2	4.7e-6
<i>I-MAD</i> (ST+)	97.5	97.9	97.7	N/A	97.7	N/A

The network includes an embedding layer, two convolutional layers with large filters and strides, a global max-pooling layer, and two fully-connected layers. The output of one convolutional layer serves as the gate of the other.

- **Krcal2018Conv** Following Raff et al. [43], Krcal et al. [32] treat an executable as a sequence of bytes and apply a CNN for malware detection, but their CNN is deeper and has smaller filters. There are four convolutional layers and four fully connected layers. Instead of a global max-pooling layer, they use a global mean-pooling layer after the convolutional layers.
- **Mour2019CNN** Mourtaji et al. [38] convert malware binaries to grayscale images and apply a 2D CNN on malware images for malware classification.

For the papers in which the authors describe multiple ways to select features and/or apply multiple machine learning models ([37, 6]), we try with all possible settings and report the best results that their methods can achieve to compare with our model.

As the ablation study, we also compare our full model "*I-MAD* (ST+)" with "*I-MAD* (no code)" and "*I-MAD* (ST)". "*I-MAD* (no code)" is our model without using assembly code. These comparisons can show the effectiveness of modeling assembly code with *Galaxy Transformer*. "*I-MAD* (ST)" is to build the *Galaxy Transformer* with the original *Star Transformer*, rather than the *Star-Plus Transformer*, to show the effectiveness of our modifications.

We also compare our model with an SVM model that uses the same features as *I-MAD* except for assembly code, since it is not a vectorial feature. We consider linear, polynomial, and RBF kernels and use grid search for tuning hyperparameters. Comparing this baseline with *I-MAD* (no code), we can separately show the effectiveness of the feature set and our model.

5.3. Experiment Settings

We evaluate the models under two different experiment settings. The main evaluation metric is accuracy (Acc), but

we also evaluate the models with precision (P), recall (R), and F1.

- **K-Fold Cross-Validation** We first evaluate our model and others with k-fold cross-validation where $k = 5$. The original dataset is randomly split into 5 even subsets. Each subset takes a turn to be chosen as the test set. Another subset takes a turn to be chosen as the validation set. The other 3 subsets form the training set. Thus, we have ${}^5P_2 = 20$ different experiment groups. Each group contains 138,000 samples in the training set, 46,000 in the validation set, and 46,000 in the test set. We acquire the experiment results of the 20 groups and report the averages.
- **Time Split Evaluation** In addition to cross-validation evaluation, we also evaluate the models in a more challenging and realistic scenario. In real life, a malware detection system is expected to detect new malware with its knowledge of known malware. To evaluate this ability of the models, we follow Saxe et al. [48] to perform a time split experiment. We use the executables compiled before 2015 to form the training and validation set, and those compiled after 2017 to form the test set. We exclude samples with a compilation time before 2000 or after 2020, either because the compilation dates are fake or the samples are outdated. There are 106,000 samples in the training set, 20,000 in the validation set, and 40,000 in the test set. We run each model with different initialization and random seeds 5 times and report the averages of the results.

5.4. Results

The results of the k-fold cross-validation and the time split experiments are shown in Table 5 and Table 6, respectively.

The full version of *I-MAD* achieves statistically significantly better accuracy and F1 than the other models in all experiments, as the p-values in t-test are much smaller than 0.01. The improvements of our model on accuracy and F1

Table 6

Results of time split experiment. It includes the p-values of t-test for F1 and accuracy between *I-MAD* (ST+) and other models.

Model	P	R	F1	pv (F1)	Acc	pv (Acc)
<i>Mosk2008OB</i>	88.6	88.6	88.6	1.2e-15	88.6	3.4e-22
<i>Bald2013Meta</i>	88.3	88.1	88.2	1.6e-17	88.2	1.2e-22
<i>Saxe2015Deep</i>	87.4	87.7	87.5	1.4e-17	87.5	2.6e-23
<i>Raff2017MaIC</i>	88.5	89.0	88.7	1.2e-16	88.7	4.5e-22
<i>Krcal2018Conv</i>	84.2	83.2	83.7	3.1e-18	83.8	1.2e-27
<i>Mour2019CNN</i>	57.0	56.6	56.8	4.3e-31	56.9	7.1e-34
<i>SVM (same features)</i>	89.2	88.8	89.0	7.3e-15	89.0	6.1e-22
<i>I-MAD</i> (no code)	89.4	89.6	89.5	3.2e-15	89.5	6.7e-21
<i>I-MAD</i> (ST)	91.1	91.1	91.1	8.6e-11	91.1	2.6e-15
<i>I-MAD</i> (ST+)	91.4	91.6	91.5	N/A	91.5	N/A

are larger in the time split experiments than in the cross-validation experiment. Even though we make sure there is no redundancy in the dataset, some pieces of malware could be extensively similar to each other if they are from the same family and compiled with slightly different modifications. Also, their compilation time is usually close to each other. In the time split experiment, the executables in the test set are compiled at least 2 years later than any executable in the training and validation sets. This is a more difficult setting that can be reflected in the pervasively lower accuracy in the time split setting than in the cross-validation setting. Thus, the significantly larger improvement of our detection model over other models in the time split experiment indicates that it has better abilities to learn robust and consistent patterns from old samples that can be generalized to classify new samples.

It is clear that with modelling assembly code with the *Galaxy Transformer*, *I-MAD* achieves much better results than it does without modelling the assembly code. This shows that modelling assembly code with our *Galaxy Transformer* helps in differentiating malicious and benign executables. We can also see that the *Galaxy Transformer* built with *Star-Plus Transformer* (*I-MAD* (ST+)) is more effective than the one built with the original *Star Transformer* (*I-MAD* (ST)). This confirms that our modifications are useful.

SVM with the same features as *I-MAD* except for assembly code, achieves accuracy similar to other best baseline models in the cross-validation experiment, and it achieves better accuracy than other baseline methods in the time split experiment, while worse than *I-MAD* (no code). This shows that the feature set we propose is effective, and our *IFFNN* has advantages in classification performance on the same feature set.

That being said, other models, except *Mour2019CNN*, also achieve reasonably good results in all experiments. However, none of the models consistently achieves the second-best performance in both experiment settings. Even though *Saxe2015Deep* uses features from multiple scopes, they do not show better results than *Bald2013Meta* and *Mosk2008OB*. The lack of any mechanism to understand assembly code is an obvious reason, as modelling assembly code in our model

improves the performance. Our improved way of representing printable string and PE import features, combined with our *IFFNN*, is the other reason. This is validated in the next subsection.

Mour2019CNN performs much worse than other models, even though we tried alternative hyper-parameter values in addition to the values the authors provided. One reason is that the way it represents an executable as an image is not sophisticated; even a small offset change in an executable would result in totally different textures in its image. In addition, we also observe overfitting, as its accuracy on the training set achieves 89.2%, while on the test set it is 71.8%. Even though our model is also a deep learning model, it does not suffer from the overfitting problem because we use two pre-training tasks to adequately train the *Satellite-Planet Transformer* and *Planet-Star Transformer* with the rich information embedded in assembly code. In contrast, *Mour2019CNN* can only be trained with the labels of executables, which is insufficient.

5.5. Interpretability

5.5.1. Case Study

Table 1 shows how our model interprets the detection result of a sample. The primary factors that lead to the prediction of 05c199.exe to be malicious and the main assembly functions related to the prediction are given. It can be seen that the assembly code of the target executable is the primary reason. The two assembly functions that contribute the most to the prediction set the program to sleep for a certain time and then download and run an embedded executable from a remote address.

5.5.2. Qualitative Analysis

To better understand the impacts of the features we use, Table 7 shows the ten most frequent main factors leading to the prediction of a sample to be malware or benign.

Main factors for both classes The assembly code of an executable is one of the most frequent factors influencing the prediction of an executable to be malicious or benign. This means that the vector representing the semantic meaning of assembly code computed by our *Galaxy Transformer*

Table 7

Most frequent main factors leading to the predictions of the malicious or benign class.

Main factors leading to the prediction of the malicious class
Assembly code
Total number of PE imports
Number of uncommon strings
The frequency of the string "Password"
The import of <i>KERNEL32.dll</i>
Total number of strings
The import of <i>WriteFile</i>
The frequency of the string "\x02\x02GetLastError"
Subsystem
Maximum entropy of sections
Main factors leading to the prediction of the benign class
Total number of strings
Number of uncommon strings
The import of <i>LCMapStringW</i>
Total number of PE imports
Assembly code
Maximum entropy of sections
The frequency of the string "\r\x01\x01\x01\x05"
The frequency of the string "\r\x01\x01\x05\x05"
The import of <i>initterm</i>
Mean entropy of sections

is very effective for malware detection. We randomly examine the assembly functions of some malware that acquire the largest attention by the relay node at the top layer of the *Star-Galaxy Transformer*. Many of them concern malicious behaviors, such as installing itself into some registry, hijacking some common legitimate DLLs, and injecting itself into another process. We find that there are statistical differences between the two classes in the mean values of total number of strings, number of uncommon strings, total number of PE imports, and maximum entropy of the sections. To be more specific, on average there are less common strings, more uncommon strings, less PE imports, and higher entropy among malware. The fact that these features could be main factors for both classes also shows the superiority of our *IFFNN* over logistic regression: as the number of uncommon strings and the number of PE imports always have non-negative values, when each of them serves as a main factor leading to the prediction of the malicious class, its weight is positive (i.e., $w(v_j)v_j > 0 \& v_j > 0 \Rightarrow w(v)_j > 0$), and when it serves as a main factor leading to the prediction of the benign class, its weight is negative (i.e., $w(v)_j v_j < 0 \& v_j > 0 \Rightarrow w(v)_j < 0$). This cannot be achieved by logistic regression because when it is trained, the weight for each feature is determined and stays static, irrelevant of the input samples. However, the weight of each feature in *IFFNN* is dynamically computed based on the whole context, i.e., the vector representing all features.

The explanation from the perspective of statistics is as follows. All supervised machine learning classification models work by identifying the correlation between a feature and

a class. Logistic regression can only learn the independent correlation between a feature and a class, without considering the correlation between features; therefore, it is linear and the weight for each feature is static. *IFFNN* learns the correlation between a feature and a class in a context considering the correlations between different features.

Main factors for malicious class The import of "KERNEL32.dll" is a main factor for the prediction of malicious class because malware relies heavily on a large number of core APIs in it to manipulate memory and the file system. The "WriteFile" function is also a main factor because malware such as ransomware and worms uses it to write content to the file system. The string of "Password" is another main factor that more frequently appears in malware created for credential theft purposes. Malware often uses mutex for different reasons. For example, it can be used as a locking mechanism to serialize access to a resource on the system or to avoid more than one instance of itself running. "GetLastError" is used to determine whether a mutex already exists. This is the reason why the frequency of string "\x02\x02GetLastError" is a main factor leading to the prediction of malware.

Main factors for benign class "LCMapStringW" is often used by benign software to convert all characters of strings to upper/lower case, which is a feature much less used in malware. "initterm" is used by core libraries to initialize a function pointer table and does not need to be imported by software programs, and therefore it is an indicator of some benign libraries. "\r\x01\x01\x01\x05" and "\r\x01\x01\x05\x05" are two strings that appear 1.8 and 3.6 times respectively more frequently among benign executables than malicious executables.

5.5.3. Quantitative Analysis

We also use a quantitative measure to analyze the interpretation of *I-MAD*. We compute the *Gini importance (GI)* and *information gain (IG)* of the features, and then rank them based on those criteria. We then rank the features by the frequencies that they serve as the main factors for the predictions. Features serving as main factors more frequently should be relatively important features for malware detection. It should be noted that even though the importance ranked this way is relevant to the rank by *Gini importance* or *information gain*, they are not supposed to be equivalent. Even if the attribution mechanism of *I-MAD* gives a perfect interpretation, the feature importance rank based on that would still be different from the rank by *Gini importance* or *information gain*.

Table 8 shows the Spearman's Rank Correlation Coefficient between the rank given by *I-MAD*, *Gini importance*, and *information gain*. It can be seen that the Spearman's Rank Correlation Coefficient between the rank given by *I-MAD* and those given by *Gini importance* and *information gain* are 0.59 and 0.55, respectively. This shows a strong correlation between them. The correlation coefficient between the rank by information gain and by *Gini importance* is only

Table 8

The Spearman's Rank Correlation Coefficient between the feature importance rank given by *I-MAD*, *Gini importance*, and *information gain*.

	IG	GI	<i>I-MAD</i>
IG	1.0	0.72	0.59
GI	0.72	1.0	0.55

Table 9

Efficiency of each model in terms of number of samples classified per second. The time consumption for feature extraction is not included.

Model	n samples per second
<i>Mosk2008OB</i>	32,152
<i>Bald2013Meta</i>	127,988
<i>Saxe2015Deep</i>	142,711
<i>Raff2017MalC</i>	86
<i>Krcal2018Conv</i>	142
<i>Mour2019CNN</i>	391
<i>SVM (same features)</i>	58
<i>I-MAD (no code)</i>	28,197
<i>I-MAD (ST)</i>	15,355
<i>I-MAD (ST+)</i>	15,239

0.72, even though they are often used for the exact same purpose: feature selection. The result means that the *IFFNN* in *I-MAD* frequently uses features that have high information gain or Gini importance as its main classification factors.

5.6. Efficiency Study

The efficiency of *I-MAD* and all models for comparison measured by the number of samples classified per second is presented in Table 9.

Among all models, the efficiency of *I-MAD* is moderate. And *I-MAD* is the second most efficient deep learning model. *Saxe2015Deep* is the most efficient because the dimension of its feature vector is only 1024, and the network is very small. *Raff2017MalC* and *Krcal2018Conv* are slow because they rely on whole byte sequences, and they are very computationally expensive. With our Titan Xp graphics cards their batch sizes could be around 32 and 128 at most, respectively. The batch size for *Mour2019CNN* depends on the number of bytes in the samples; in extreme cases we need to run the model on the CPU because the graphics card memory cannot hold the computation for even one large executable. For *I-MAD (ST)/(ST+)*, the batch size could be at least 512 for most samples. As the representation of assembly code is computed at three levels (i.e., basic block, function, and executable), the memory for the lower level computation is released and reused after the representation is computed. For *I-MAD (no code)*, the batch size could be 5,120. It is worth the extra computational cost to model assembly code because the benefit of it in classification performance is significant. *SVM* with the same features as *I-MAD* also has very low efficiency because its computational com-

plexity is linear with the dimension of feature vector and the number of support vectors, which are large when the dataset is complex. In our experiments, there are always more than 43,000 support vectors, and the dimension of feature vectors is more than 2,700.

6. Limitations and Future Work

Adversarial attack and defense are closely related topics to classification problems such as image classification [26] and malware detection [4]. Interpretability is a double-edged sword considering adversarial attacks in white-box settings, where adversaries have full access to the *I-MAD* model and can use the interpretations to craft adversarial samples more easily than by using an uninterpretable model.

Evasion techniques (e.g., adversarial attacks) are extensively applied in wild malware, as is the case of our dataset. Following previous experience [3, 29], we counter the evasion techniques by detecting malware from the views of multiple disparate feature sets. Also, since adversarial samples are already in our dataset, adversarial training is automatically performed to defend against the attacks [26]. Usually, adversarial attack and defense are discussed in a different research paper than the one proposing a novel classification model. One of our future work directions is to further investigate adversarial attack and defense on malware detection.

7. Conclusion

In this paper, we present *I-MAD*, a novel deep learning model for static malware detection that is based on the understanding of assembly code. In addition to its excellent detection performance, it can also provide interpretation for its detection results, which can be examined by malware analysts. Therefore, in addition to malware detection, it can also help malware analysts locate malicious payloads and find consistent patterns in malware samples. The proposed *IFFNN* has values that can be applied in interpretable classification for other tasks as well.

CRedit authorship contribution statement

Miles Q. Li: Conceptualization, Methodology, Software, Validation, Data Curation, Investigation, Writing - Original Draft. **Benjamin C. M. Fung:** Conceptualization, Validation, Supervision, Writing - Review & Editing, Funding acquisition. **Philippe Charland:** Project administration, Writing - Review & Editing. **Steven H. H. Ding:** Data Curation, Writing - Review & Editing.

Acknowledgments

This research is supported by Defence Research and Development Canada (contract no. W7701-176483/001/QCL), NSERC Discovery Grants (RGPIN-2018-03872), and Canada Research Chairs Program (950-232791). The Titan Xp used for this research was donated by the NVIDIA Corporation.

References

- [1] Amer, E., Zelinka, I., 2020. A dynamic windows malware detection and prediction method based on contextual understanding of api call sequence. *Computers & Security* 92, 101760.
- [2] Anderson, B., Quist, D., Neil, J., Storlie, C., Lane, T., 2011. Graph-based malware detection using dynamic analysis. *Journal in Computer Virology* 7, 247–258.
- [3] Anderson, B., Storlie, C., Lane, T., 2012. Improving malware classification: bridging the static/dynamic gap, in: *Proceedings of the 5th ACM workshop on Security and artificial intelligence*, ACM. pp. 3–14.
- [4] Anderson, H.S., Kharkar, A., Filar, B., Evans, D., Roth, P., 2018. Learning to evade static pe machine learning malware models via reinforcement learning. *arXiv preprint arXiv:1801.08917*.
- [5] Bahdanau, D., Cho, K., Bengio, Y., 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- [6] Baldangombo, U., Jambaljav, N., Horng, S.J., 2013. A static malware detection system using data mining methods. *arXiv preprint arXiv:1308.2831*.
- [7] Basole, S., Di Troia, F., Stamp, M., 2020. Multifamily malware models. *Journal of Computer Virology and Hacking Techniques*, 1–14.
- [8] Bayer, U., Moser, A., Kruegel, C., Kirda, E., 2006. Dynamic analysis of malicious code. *Journal in Computer Virology* 2, 67–77.
- [9] Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- [10] Caruana, R., Lawrence, S., Giles, C.L., 2001. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping, in: *Advances in Neural Information Processing Systems*, pp. 402–408.
- [11] Cerna, A.E.U., Pattichis, M., VanMaanen, D.P., Jing, L., Patel, A.A., Stough, J.V., Haggerty, C.M., Fornwalt, B.K., 2019. Interpretable neural networks for predicting mortality risk using multi-modal electronic health records. *arXiv preprint arXiv:1901.08125*.
- [12] Cesare, S., Xiang, Y., 2010. Classification of malware using structured control flow, in: *Proceedings of the Eighth Australasian Symposium on Parallel and Distributed Computing*, Australian Computer Society, Inc.. pp. 61–70.
- [13] Cesare, S., Xiang, Y., Zhou, W., 2013. Control flow-based malware variant detection. *IEEE Transactions on Dependable and Secure Computing* 11, 307–317.
- [14] Chen, J., Alalfi, M.H., Dean, T.R., Zou, Y., 2015. Detecting android malware using clone detection. *Journal of Computer Science and Technology* 30, 942–956.
- [15] Choi, E., Bahadori, M.T., Sun, J., Kulas, J., Schuetz, A., Stewart, W., 2016. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism, in: *Advances in Neural Information Processing Systems*, pp. 3504–3512.
- [16] Dahl, G.E., Stokes, J.W., Deng, L., Yu, D., 2013. Large-scale malware classification using random projections and neural networks, in: *Proceedings of the 38th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE. pp. 3422–3426.
- [17] Dai, J., Guha, R.K., Lee, J., 2009. Efficient virus detection using dynamic instruction sequences. *JCP* 4, 405–414.
- [18] Damodaran, A., Di Troia, F., Visaggio, C.A., Austin, T.H., Stamp, M., 2017. A comparison of static, dynamic, and hybrid analysis for malware detection. *Journal of Computer Virology and Hacking Techniques* 13, 1–12.
- [19] Demontis, A., Melis, M., Biggio, B., Maiorca, D., Arp, D., Rieck, K., Corona, I., Giacinto, G., Roli, F., 2017. Yes, machine learning can be more secure! a case study on android malware detection. *IEEE Transactions on Dependable and Secure Computing*.
- [20] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [21] Ding, S.H., Fung, B.C., Charland, P., 2019. Asm2vec: Boosting static representation robustness for binary clone search against code obfuscation and compiler optimization, in: *IEEE Symposium on Security and Privacy (SP)*, IEEE. pp. 472–489.
- [22] Dong, L., Yang, N., Wang, W., Wei, F., Liu, X., Wang, Y., Gao, J., Zhou, M., Hon, H.W., 2019. Unified language model pre-training for natural language understanding and generation, in: *Advances in Neural Information Processing Systems*, pp. 13063–13075.
- [23] Dong, Y., Su, H., Zhu, J., Bao, F., 2017. Towards interpretable deep neural networks by leveraging adversarial examples. *arXiv preprint arXiv:1708.05493*.
- [24] Fredrikson, M., Jha, S., Christodorescu, M., Sailer, R., Yan, X., 2010. Synthesizing near-optimal malware specifications from suspicious behaviors, in: *IEEE Symposium on Security and Privacy (SP)*, IEEE. pp. 45–60.
- [25] Gibert, D., Mateu, C., Planes, J., 2020. Hydra: A multimodal deep learning framework for malware classification. *Computers & Security*, 101873.
- [26] Goodfellow, I.J., Shlens, J., Szegedy, C., 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- [27] Guo, Q., Qiu, X., Liu, P., Shao, Y., Xue, X., Zhang, Z., 2019. Star-transformer. *arXiv preprint arXiv:1902.09113*.
- [28] Huang, W., Stokes, J.W., 2016. Mtnet: a multi-task neural network for dynamic malware classification, in: *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*, Springer. pp. 399–418.
- [29] Islam, R., Tian, R., Batten, L.M., Versteeg, S., 2013. Classification of malware based on integrated static and dynamic features. *Journal of Network and Computer Applications* 36, 646–656.
- [30] Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [31] Kolter, J.Z., Maloof, M.A., 2004. Learning to detect malicious executables in the wild, in: *Proceedings of the tenth ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, ACM. pp. 470–478.
- [32] Krčál, M., Švec, O., Bálek, M., Jašek, O., 2018. Deep convolutional malware classifiers can learn from raw executables and labels only. *International Conference on Learning Representations (ICLR) 2018 Workshop*.
- [33] Kruegel, C., Kirda, E., Mutz, D., Robertson, W., Vigna, G., 2005. Polymorphic worm detection using structural information of executables, in: *International Workshop on Recent Advances in Intrusion Detection*, Springer. pp. 207–226.
- [34] Lipton, Z.C., 2018. The mythos of model interpretability. *Queue* 16, 31–57.
- [35] Lou, Y., Caruana, R., Gehrke, J., 2012. Intelligible models for classification and regression, in: *Proceedings of the 18th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pp. 150–158.
- [36] Luong, M.T., Pham, H., Manning, C.D., 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- [37] Moskovitch, R., Feher, C., Tzachar, N., Berger, E., Gitelman, M., Dolev, S., Elovici, Y., 2008. Unknown malcode detection using opcode representation, in: *Intelligence and Security Informatics*. Springer, pp. 204–215.
- [38] Mourtaji, Y., Bouhorma, M., Alghazzawi, D., 2019. Intelligent framework for malware detection with convolutional neural network, in: *Proceedings of the 2nd International Conference on Networking, Information Systems & Security*, ACM. p. 7.
- [39] Onwuzurike, L., Mariconti, E., Andriotis, P., Cristofaro, E.D., Ross, G., Stringhini, G., 2019. Mamadroid: Detecting android malware by building markov chains of behavioral models (extended version). *ACM Transactions on Privacy and Security (TOPS)* 22, 1–34.
- [40] Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A., 2017. Automatic differentiation in pytorch. *Neural Information Processing Systems NIPS 2017 Autodiff Workshop*.
- [41] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., 2018. Improving language understanding with unsupervised learning. *Techni-*

- cal Report. Technical report, OpenAI.
- [42] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., 2019. Language models are unsupervised multitask learners. OpenAI Blog 1, 8.
- [43] Raff, E., Barker, J., Sylvester, J., Brandon, R., Catanzaro, B., Nicholas, C., 2017. Malware detection by eating a whole exe. arXiv preprint arXiv:1710.09435 .
- [44] Raff, E., Zak, R., Cox, R., Sylvester, J., Yacci, P., Ward, R., Tracy, A., McLean, M., Nicholas, C., 2018. An investigation of byte n-gram features for malware classification. Journal of Computer Virology and Hacking Techniques 14, 1–20.
- [45] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. arXiv preprint arXiv:1910.10683 .
- [46] Royal, P., Halpin, M., Dagon, D., Edmonds, R., Lee, W., 2006. Polyunpack: Automating the hidden-code extraction of unpack-executing malware, in: Proceedings of the 22nd Annual Computer Security Applications Conference (ACSAC'06), IEEE. pp. 289–300.
- [47] Santos, I., Devesa, J., Brezo, F., Nieves, J., Bringas, P.G., 2013. Opem: A static-dynamic approach for machine-learning-based malware detection, in: Proceedings of the International Joint Conference CISIS'12-ICEUTE 12-SOCO 12 Special Sessions, Springer. pp. 271–280.
- [48] Saxe, J., Berlin, K., 2015. Deep neural network based malware detection using two dimensional binary program features, in: Proceedings of the 10th International Conference on Malicious and Unwanted Software (MALWARE), IEEE. pp. 11–20.
- [49] Schultz, M.G., Eskin, E., Zadok, F., Stolfo, S.J., 2001. Data mining methods for detection of new malicious executables, in: Proceedings 2001 IEEE Symposium on Security and Privacy. S&P 2001, IEEE. pp. 38–49.
- [50] Shickel, B., Tighe, P.J., Bihorac, A., Rashidi, P., 2017. Deep ehr: a survey of recent advances in deep learning techniques for electronic health record (ehr) analysis. IEEE Journal of Biomedical and Health Informatics 22, 1589–1604.
- [51] Vasan, D., Alazab, M., Wassan, S., Safaei, B., Zheng, Q., 2020. Image-based malware classification using ensemble of cnn architectures (imcec). Computers & Security , 101748.
- [52] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need, in: Advances in Neural Information Processing Systems, pp. 5998–6008.
- [53] Vemparala, S., Di Troia, F., Corrado, V.A., Austin, T.H., Stamo, M., 2016. Malware detection using dynamic birthmarks, in: Proceedings of the 2016 ACM on international workshop on security and privacy analytics, pp. 41–46.
- [54] Verma, V., Mutttoo, S.K., Singh, V., 2020. Multiclass malware classification via first-and second-order texture statistics. Computers & Security 97, 101895.
- [55] Wong, W., Stamp, M., 2006. Hunting for metamorphic engines. Journal in Computer Virology 2, 211–229.
- [56] Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., Le, Q.V., 2019. Xlnet: Generalized autoregressive pretraining for language understanding. arXiv preprint arXiv:1906.08237 .
- [57] Ye, Y., Li, T., Adjeroh, D., Iyengar, S.S., 2017. A survey on malware detection using data mining techniques. ACM Computing Surveys (CSUR) 50, 41.
- [58] Zeiler, M.D., Fergus, R., 2014. Visualizing and understanding convolutional networks, in: European conference on computer vision, Springer. pp. 818–833.

Miles Q. Li is a Ph.D. candidate in the School of Computer Science, McGill University, Montreal, Canada. He received his B.Sc and M.Sc from Peking University. His research interests include data mining for cybersecurity, deep learning, and natural language processing.

Benjamin C. M. Fung is a Canada Research Chair in Data Mining for Cybersecurity, a Full Professor with the School of Information Studies, and an Associate Member with the School of Computer Science at McGill University in Canada. He received a Ph.D. degree in computing science from

Simon Fraser University in Canada in 2007. He has over 130 refereed publications that span the research forums of data mining, privacy protection, cybersecurity, services computing, and building engineering. His data mining works in crime investigation and authorship analysis have been reported by media worldwide. Prof. Fung is a licensed Professional Engineer of software engineering in Ontario, Canada.

Philippe Charland is a Defence Scientist at Defence Research and Development Canada – Valcartier Research Centre, in the Mission Critical Cyber Security Section. His research interests encompass software reverse engineering and computer forensics. As a member of the Systems Vulnerabilities and Lethality Group, his research focuses on binary-level program comprehension to accelerate the reverse engineering process involved in malicious software analysis and classification, as well as for mission assurance. Mr. Charland holds a bachelor and a master's degree in Computer Science, both from Concordia University, Montreal, Canada.

Steven H. H. Ding is an Assistant Professor in the School of Computing at Queen's University, where he leads the LINNA Artificial Intelligence and Security Lab. His research bridges the domain of machine learning, data mining, and cybersecurity, aiming at addressing critical cybersecurity challenges using AI technologies and securing the future of AI systems. Dr. Ding obtained his Ph.D. from McGill University in 2019, and he was awarded the FRQNT Doctoral Research Scholarship of Quebec and the Dean's Graduate Award at McGill University.