

# MALIGN: Adversarially Robust Malware Family Detection using Sequence Alignment

Shoumik Saha  
Department of CSE, BUET  
Dhaka, Bangladesh  
shoumiksaha901@gmail.com

Sadia Afroz  
Avast, ICSI  
sadia@icsi.berkeley.edu

Atif Rahman  
Department of CSE, BUET  
Dhaka, Bangladesh  
atif@cse.buet.ac.bd

**Abstract**—We propose MALIGN, a novel malware family detection approach inspired by genome sequence alignment. MALIGN encodes malware using four nucleotides and then uses genome sequence alignment approaches to create a signature of a malware family based on the code fragments conserved in the family making it robust to evasion by modification and addition of content. Moreover, unlike previous approaches based on sequence alignment, our method uses a multiple whole-genome alignment tool that protects against adversarial attacks such as shuffling of code blocks. Our approach outperforms state-of-the-art machine learning based malware detectors and demonstrates robustness against trivial adversarial attacks. MALIGN also helps identify the techniques malware authors use to evade detection.

**Index Terms**—malware, adversarial, sequence alignment

## I. INTRODUCTION

To detect the rising number of malwares at scale, machine learning is necessary. Indeed, currently all commercial malware detectors use machine learning. However, one shortcoming of current static malware detectors is that they can be easily evaded by changing a malware trivially without changing the core of the malware [1]–[10]. Fundamentally, the adversarial attacks use one simple technique: add or modify selected content to the malware.

One simple way to design an adversarially robust malware detector is to make it rely on the core functionalities of a malware and ignore the added insignificant content. This turns out to be surprisingly hard, especially for static detectors, without increasing the false positive rate. This happens because malware authors often take a similar approach to evade detection—hide malicious content in the overlay and data section instead of the code section without changing the malware semantics.

Taking inspiration from bioinformatics, we model a malware like a DNA sequence or a genome. Just as DNA sequences are made of only four types of nucleotides, malwares are sequences of bits, and modifications of malwares mirror accumulation of mutations in genomes during evolution. Genomes contain critical regions for the survival of the organism, such as, protein coding genes where mutations may be lethal. Similarly, malwares contain code blocks that are difficult to modify without altering its functionality and semantics. If we can translate a malware in terms of the basic

building blocks, our detector will be robust by design that cannot be evaded without substantially changing the malware.

We propose MALIGN, a novel malware family detection approach inspired by genome sequence alignment. Our approach at first converts a family of malware files into malware nucleotide sequence files i.e. sequences of A, C, G and T. Then we use a multiple whole-genome alignment tool to identify common alignment blocks per family. These alignment blocks work as a signature of the malware family, and a score is assigned to each of them depending on their importance. We train a classifier using the features that represent how well a block identifies with a particular family. To classify whether a new malware belongs to the family, we first compute the alignment of the new malware with the sequences representing the blocks i.e. signature of the family, and use it to classify the malware.

Our robustness properties come from the use a recent multiple whole-genome alignment method that can find conserved blocks of sequences even in the presence of sequence re-ordering and minor modifications, and through estimation of the degree of conservation at each location by processing the generated alignment. It prevents certain types of adversarial manipulation, such as, adding extra content, changing code order, and minor changes to the code. To evade detection, an attacker needs to change the code significantly.

We evaluate MALIGN on two datasets: Kaggle Microsoft Malware Classification Challenge (Big 2015) and Microsoft Machine Learning Security Evasion Competition (2020) (MLSec). In comparison with MalConv, feature fusion and CNN-based malware classifiers, our approach has higher accuracy and robustness. Moreover, sequence alignment helps to reveal the common practices of malware families, such as hiding code in a non-code sections.

In summary, our main contributions are:

- **Scalable and Explainable:** Our approach is simple, scalable and easily explainable due to the use of new multiple sequence alignment tool.
- **High accuracy with low train data:** MALIGN achieves high accuracy even without large training data.
- **Robustness to adversarial attacks:** MALIGN finds conserved code blocks considering mismatches and assigns high scores to those. So, critical code blocks may need to be modified drastically to evade detection.

## II. BACKGROUND

Sequence alignment is a widely studied problem in bioinformatics to find similarity among DNA, RNA or protein sequences, and to study evolutionary relationships among diverse species. It is the process of arranging sequences in such a way that regions of similarity are *aligned*, with gaps (denoted by '-') inserted to represent insertions and deletions in sequences. An alignment of the sequences ATTGACCTGA and ATCGTGTA is shown below where the regions denoted in black, characterized by matched characters, are *conserved* whereas the red and blue regions denote substitutions i.e. point mutations, and insertions or deletions during the evolutionary process respectively.

```
ATTGACCTG-A
AT---CGTGA
```

In sequence alignment, matches, mismatches and insertions-deletions (in-dels) are assigned scores based on their frequencies during evolution and the goal is to find an alignment with the maximum score. The problem of finding an optimal alignment of the entire sequences (global alignment) and that of finding an optimal alignment of their sub-sequences (local alignment) can be solved by dynamic programming using the Needleman-Wunsch [11] and Smith-Waterman [12] algorithms respectively. While the algorithms can be used to align more than two sequences, the running time is exponential in the number of sequences. To address the tractability issue, a number of tools have been developed [13]–[15], that use heuristics to solve the multiple sequence alignment (MSA) problem.

However, in addition to point mutations and short insertions-deletions, large scale genome rearrangement events take place during evolution. Such genome rearrangement events include reversal of a genomic segment (inversion), shuffling of order of genomic segments (transposition or translocation), duplication and deletion of segments, etc. Although the aforementioned tools are unable to deal with genome rearrangements, methods such as MUMmer [16] can perform alignment of two sequences in presence of rearrangements whereas Mauve [17], Cactus [18], etc. can handle multiple sequences.

Recently, Armstrong *et al.* have developed Progressive Cactus [19], and Minkin & Medvedev have developed SibeliaZ [20] that can align hundreds to thousands of whole-genome sequences in presence of rearrangements. The tools identify similar sub-sequences in the sequences from different species to create blocks of rearrangement-free sequences, and then performs a multiple sequence alignment of the sequences in each block.

Since adversaries can modify malwares relatively easily by changing orders of blocks of codes without altering functionality of the malware, it is important that the tool used to align malware sequences is robust to such rearrangements in code. Here, we use SibeliaZ to align malware sequences to identify conserved blocks of codes and calculate a conservation score of the blocks for malware detection and classification. It is

worth noting that the blocks of codes identified need not be fully conserved, i.e. there can be modifications, insertions and deletions of small number of instructions within the blocks, making it robust to adversarial attacks.

## III. RELATED WORK

To counter the increasing amount of malware and detect them, several methods and techniques have been developed over the years. In the early days, Wressnegger *et al.* [21] and Zakeri *et al.* [22] proposed a signature based approach using static analysis. Later, a dynamic approach - malware detection by analysing the malware behavior, was proposed by Martignoni *et al.* [23] and Willems *et al.* [24]. In recent times, machine learning based techniques are mostly being used to classify malwares. Schultz *et al.* [25] first proposed a data mining technique for malware detection using three different types of static features. Subsequently, Nataraj *et al.* [26] proposed a malware classification approach based on image processing techniques by converting the bytes files to image files. Later, Kalash *et al.* [27] improved on [26] by developing M-CNN using malware images. Besides CNN, RNN has also been used for malware analysis. [28] and [29] proposed techniques with LSTM using opcode sequences of malware. Santos *et al.* [30] proposed a hybrid technique by integrating both static and dynamic analysis. Yan *et al.* [31] developed MalNet using an ensemble method on CNN, LSTM and extracting metadata features. Ahmadi *et al.* [32] extracted and selected features of malware depending on the importance and applied feature fusion on them. Recently, Raff *et al.* [33] developed a state-of-the-art technique, MalConv using only the raw byte sequence as the input to a neural network.

Prior work proposed two main ways to improve the adversarial robustness of malware detectors: adversarial training, and robustness by design. Adversarial training, where a malware detector is trained with adversarial examples, is one of the mostly used approaches to improve adversarial robustness [34]. In the malware domain, several work demonstrated that adversarial training can improve the robustness significantly without reducing the accuracy on the original sample [35]–[37]. Robustness by design approaches build classifiers to eliminate a certain classes of adversarial attacks. Certified or provable robustness is a robustness by design approach that trains classifiers with local robustness properties that can provably eliminate classes of evasion attacks [38]. In the malware domain, Chen *et al.* [39] proposed learning PDF malware detectors with verifiable robustness properties. Íñigo *et al.* [40] trained a XGBoost based malware detector with the monotonicity property that ensures that an adversary cannot decrease the classification score by adding extra content. The first method relies on the availability of enough adversarial samples, which may not always be the case in the fast-changing malware world. Our approach falls under the second category, robustness by design. Although our approach does not provide any provable robustness guarantees, it increases the cost of an adversary by eliminating the possibility of trivial attacks.

In the past, sequence alignment based approaches have been used for malware analysis by a number of researchers [41]–[47]. Chen et al. used multiple sequence alignment to align computer viral and worm codes of variable lengths to identify invariant regions [41]. This approach was subsequently enhanced in [42], [43], [45]. However, none of these methods address the issue that blocks of code can be shuffled without affecting malware behaviour.

Sequence alignment has also been applied on system call sequences of malwares to extract evasion signatures and cluster samples [44], classify malware families [46], and for malware detection, classification and visualization [47]. While it is more difficult for malware developers to shuffle API calls without changing the behaviour of malwares, these approaches require access to API call sequences of malwares and are not suitable in all circumstances.

Drew et al. [48] utilized another approach developed by the bioinformatics and computational biology community for malware classification - that for gene or sequence classification. The method is based on extracting short words i.e.  $k$ -mers from sequences and calculating similarity between sequences based the set of words present in them. Although the method is efficient, it does not fully utilize the information provided by long stretches of conserved regions in malwares, and is not suitable for identifying critical code blocks in malwares. Moreover, MALIGN has higher accuracy rate than this method

#### IV. METHODS

##### Overview

To classify or detect known/unknown malwares and its variants, in this paper, we propose a malware classification or detection system based on multiple whole-genome alignment. The basic building block of the method is a binary classification system that can predict whether an instance belongs to a particular malware family or not. The input to this binary classifier is a training set consisting of positive examples i.e. malwares from a particular family, and negative examples which may be non-malwares or malwares from other families. The system can be extended to malware detection by creating a binary classifier for all malware families. The instances that are predicted to be negative by all these classifiers can then be treated as benign.

The main steps of our proposed method are shown in Algorithm 1.

The method is illustrated in Figure 1. We start with the given *malware bytes files* i.e. executable files and convert them into *malware nucleotide sequence files* i.e. sequences of A, C, G and T. Then these nucleotide sequence files are aligned using a multiple whole-genome alignment tool (SibeliaZ) which outputs alignment blocks that are common among these files. These alignment blocks are merged and thus *consensus sequence* is constructed. In this step, *conservation score* for each coordinate of the consensus sequence is also generated. This consensus sequence is aligned with each sample from a balanced train set with positive and negative samples with respect to the malware family of interest, and an *alignment score*

#### Algorithm 1: MALIGN

---

**Input:** Training set,  $\mathbb{X} = (\mathbb{X}_+, \mathbb{X}_-)$ , where  $\mathbb{X}_+$ : byte files from malwares of a family,  $\mathbb{X}_-$ : byte files from non-malwares or malwares from other families, and Test set,  $\mathbb{Y}$

**Output:** Labels for  $\mathbb{Y}$

$(\mathbb{S}_+, \mathbb{S}_-) \leftarrow$  Convert to nucleotide sequence files  $(\mathbb{X}_+, \mathbb{X}_-)$

$\mathbb{B} \leftarrow$  Perform multiple sequence alignment and identify conserved blocks  $(\mathbb{S}_+)$

**forall** blocks  $B \in \mathbb{B}$  **do**

$C_B \leftarrow$  Get consensus sequence  $(B)$

**for**  $i = 1 \rightarrow \text{length}(B)$  **do**

Calculate  $\text{ConservationScore}(B, i, N)$  where  $N = A, C, G, T$

**forall** sequences  $Z \in (\mathbb{S}_+ \cup \mathbb{S}_-)$  **do**

$F_Z :=$  Feature vector of  $Z$

**forall** blocks  $B \in \mathbb{B}$  **do**

$S \leftarrow$  Get alignments  $(Z, C_B)$

$\text{AlignmentScore} \leftarrow$  Calculate alignment score  $(S, B)$

$\text{AlignmentCount} \leftarrow$  Get alignment count  $(S)$

$F_Z \leftarrow F_Z \cup (\text{AlignmentScore}, \text{AlignmentCount})$

$\mathcal{M} \leftarrow$  Learn classification model  $(\mathbb{F}, \mathbb{X})$  where  $\mathbb{F}$ : feature matrix

$\mathbb{T} \leftarrow$  Convert to nucleotide sequence files  $(\mathbb{Y})$

$\mathbb{L} :$  labels

**forall** sequences  $T \in \mathbb{T}$  **do**

$F_T :=$  Feature vector of  $T$

**forall** blocks  $B \in \mathbb{B}$  **do**

$S \leftarrow$  Get alignments  $(T, C_B)$

$\text{AlignmentScore} \leftarrow$  Calculate alignment score  $(S, B)$

$\text{AlignmentCount} \leftarrow$  Get alignment count  $(S)$

$F_T \leftarrow F_T \cup (\text{AlignmentScore}, \text{AlignmentCount})$

$L \leftarrow$  Predict  $(\mathcal{M}, F_T)$

$\mathbb{L} \leftarrow \mathbb{L} \cup L$

**return**  $\mathbb{L}$

---

is calculated for every sample for each conserved block. These scores are then used as input to a machine learning model which learns a classifier to distinguish between malwares belonging to the family, and malwares from other families as well as non-malwares. To classify a new sample, the sequence is aligned with the consensus sequence and alignment scores for the new instance are generated similarly, and the scores are passed into our classifier to classify the new sample. Each of these steps is described in more details below.

##### Bytes file to nucleotide sequence file conversion

First the binary executable or bytes files are converted to nucleotide sequence files containing sequences of A, C, G and T. The conversion is performed so that existing whole-genome alignment tools can be used. The conversion from the byte code to nucleotide sequence is done by converting each pair of bits to a nucleotide according to Table I.

Byte Character	Nucleotide
00	A
01	C
10	G
11	T

TABLE I  
BYTES TO NUCLEOTIDE MAPPING

In some malware datasets such as the Kaggle Microsoft Malware Classification Challenge (Big 2015) dataset [49],

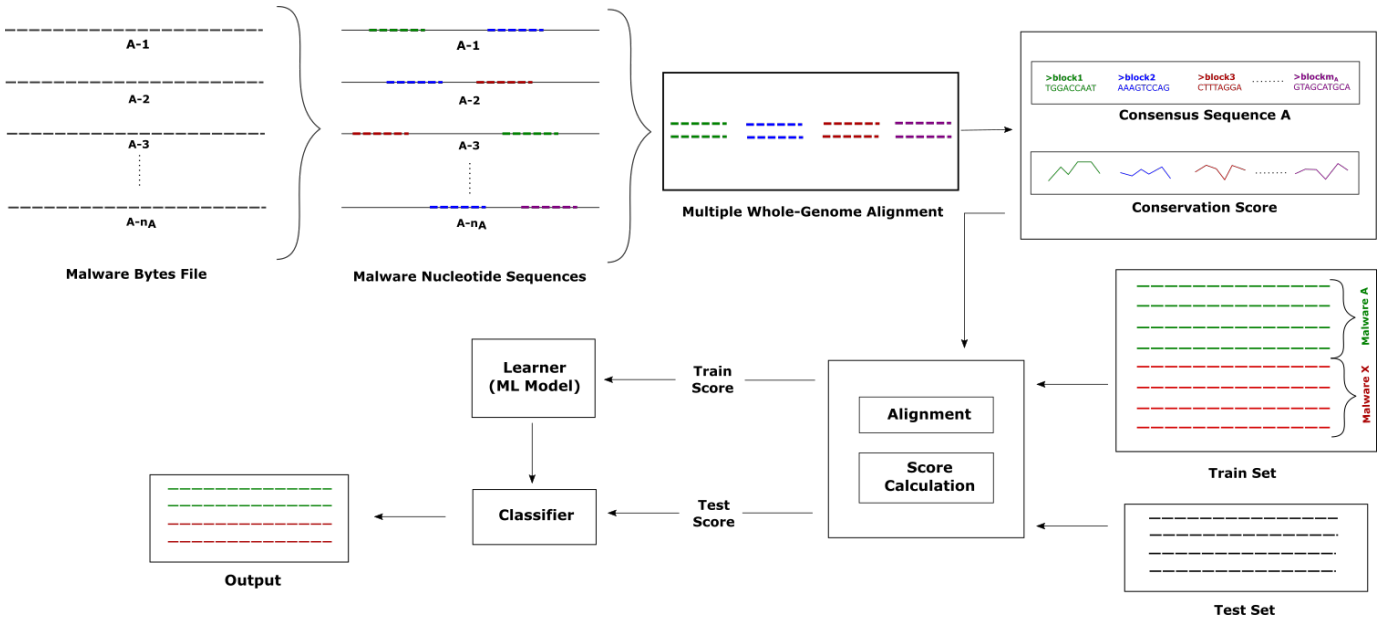


Fig. 1. Overview of MALIGN. (1) The malware bytes files (executables) from malwares of a particular family are first converted to nucleotide sequence files. (2) Then the malware nucleotide sequences are aligned using a multiple sequence alignment tool Sibeliaz. It first identifies similar sequences in different malwares to form blocks. Highly similar sequences (colored sequences) can be in different order in different files. The sequences in each block are then aligned. (3) The aligned sequences in each block are used to construct consensus sequences and conservation scores are calculated for each conserved block. (4) Then two sets of sequences - one corresponding to the malware family of interest and the other corresponding to non-malwares or malwares from other families are aligned to the consensus sequences and the degrees of conservation of each conserved block in the training sequences are estimated. (5) Finally a machine learning model is learnt to classify sequences based on the alignment scores of the sequences with the blocks. To classify new instances, sequences are aligned to the consensus sequences of the blocks and alignments are scored. The scores are then used as features for the class prediction.

the provided bytes files contain “??” and long stretches of “00” in some cases which do not preserve any significant value or meaning. These are removed before the conversion to nucleotide sequences.

#### Multiple alignment of malware nucleotide sequences

The next step is to align the malware nucleotide sequences. In this paper, we use the multiple whole-genome alignment tool Sibeliaz [20]. Sibeliaz performs whole-genome alignment of multiple sequences and constructs locally co-linear blocks. Figure 2 illustrates alignment of three different malware nucleotide sequence files from the same family. The sequences share blocks of similar sequences showed in dashed lines of same color. They may also contain sequences unique to each sequence indicated by lines with different colors.

During the block construction process:

- The order of the shared blocks may differ in different sequences and the blocks may not be shared across all sequences. This helps MALIGN to be robust to the evasion attempts, such as shuffling blocks of code.
- The shared blocks may not be fully conserved i.e. there may be mismatches of characters to some extent which means minor alteration, modification to the code will not prevent detection of blocks.

These properties of multiple whole-genome alignment bolsters the robustness of MALIGN against many obfuscation techniques.

Sibeliaz first identifies the shared linear blocks and then performs multiple sequence alignment of locally co-linear blocks. The block coordinates are output in GFF format and the alignment is in MAF format. The GFF format is a file format used for describing genes and other features of DNA, RNA and protein sequences. The multiple alignment format (MAF) is a format for describing multiple alignments in a way that is easy to parse and read. In our case, this format stores multiple alignment blocks at the byte code level among malwares. We generate such an MAF file for each malware family using training samples and identify the blocks of codes that are highly conserved across the malware family.

#### Consensus sequence and score generation

We process all sequences of the alignment blocks of MAF file from the previous step and generate a new sequence for each block, which is known as *consensus sequence* [50]. At first, we scan the length of all sequences and find the maximum one that will be the length of our consensus sequence. Then we traverse through the coordinates of every sequence and find the nucleotide of highest occurrence for each coordinate. We put the most frequently occurring nucleotide in corresponding index of the consensus sequence.

In Figure 2, consensus sequence generation of alignment block for Block-1 is shown in detail. Below the alignment block, the corresponding sequence logo is shown. The height of the individual letters in sequence logo represents how com-

mon the corresponding letter is at that particular coordinate of the alignment.

We thus construct the consensus sequence by taking the letter (nucleotide) with highest frequency for each coordinate. Similarly, the consensus sequences for all blocks are generated and are stored in a file in FASTA format with a unique id. These consensus sequences are the conserved part of the malware family which can be considered as the signature or common pattern of that family. The files are used in subsequent steps to classify malwares.

In addition to the consensus sequence, we calculate conservation scores for the blocks. In bioinformatics, conservation score is used during evaluation of sites in a multiple sequence alignment, in order to identify residues critical for structure or function. This is calculated per base, indicating how many species in a given multiple alignment match at each locus. In malware world, the conservation score can indicate the significance or importance of a code segment, or function call, or API call in a malware family. The responsible code segments of a malware will have high conservation scores compared to the segments those are not frequent, or conserved in malware files.

In Figure 2, the height of the bars of conservation score indicates the degree of conservation at the corresponding position. For each coordinate, we store the score for each of the four nucleotides which is given by the occurrence percentage of that nucleotide at that coordinate. So, at the  $i^{th}$  index of the alignment block  $B$ , the score for nucleotide  $N(= A, C, G, T)$  is given by:

$$ConservationScore_{B, i, N} = \frac{\# \text{ of occurrences of } N \text{ at the } i^{th} \text{ index in block } B}{\# \text{ of sequences in alignment block } B}$$

For example, in Figure 2, Block-1 has 3 sequences in total. Since at the 1<sup>st</sup> index the block contains 3 Gs,

$$Score_{1,1,G} = 3/3 = 1.00$$

Again at the 6<sup>th</sup> index the block contains 2 Ts and 1 A. So

$$Score_{1,6,T} = 2/3 = 0.66$$

$$Score_{1,6,A} = 1/3 = 0.33$$

#### Alignment with consensus sequences

Once the consensus sequence and the conservation scores are generated, we take a training set for each malware family. In the training set, the positive examples are samples from that malware family and the negative examples are non-malwares or malwares from other families. All samples from the training set are aligned to the consensus sequences of the corresponding family to get the aligned blocks for each sample. Using the previously generated conservation score, we calculate new scores called alignment scores for each block for all samples which will be used as features.

In Figure 1 the green and red lines indicate the positive and negative samples in the training set respectively. These samples are then aligned with the consensus sequences using the alignment tool, SibeliaZ which outputs alignments for each sample. An example of alignment score calculation is shown in Figure 2. Malware X-1 and X-2 are positive and negative sample respectively. X-1 has three aligned sequences with the consensus sequence (shown in purple) whereas X-2 has only one. Sum of scores for all aligned sequences will be the score for the corresponding block of that sample. As an example, for total score of sample X-1, we sum the scores of 3 aligned sequences. Each aligned sequence's score is the sum of the score of all coordinates.

The aligned sequence score is then multiplied by the number of sequences that constructed the corresponding block since the higher the number of sequences that generated the block, the more conserved the sequence is across the instances from that family. In Figure 2, adding all coordinate's score of sample X-1's first aligned sequence, we get 7.32. Since the corresponding consensus sequence was generated from 3 sequences, the final score for first aligned sequence will be  $21.96(= 7.32 \times 3)$ . Finally, the total alignment score for the block was calculated by adding the scores of all 3 aligned sequences.

In general, the total alignment score of a sample  $Z$  for consensus sequence  $C_B$  from the alignment block  $B$  is given by

$$AlignmentScore_{Z, C_B} = \sum_{s \in S} \left( \sum_{i=1}^{length(s)} ConservationScore_{B, j, s_i} \right) \times \text{number of sequences in block } B$$

where,  $S$  is the set of sequences from sample  $Z$  that got aligned with  $C_B$ ,  $s_i$  is the  $i$ -th nucleotide of the sequence  $s$  and  $j$  is the index of  $C_B$  where  $s_i$  was aligned.

Along with this score, we also store the total number of times the consensus sequence of a block gets aligned with the sample. In Figure 2, the consensus sequence was aligned with sample X-1 three times. Both the number of occurrences and the total alignment score for the consensus sequence of each block for a malware family are used as features for the subsequent classification, resulting in  $2m$  features if multiple sequence alignment of a malware family has  $m$  aligned blocks.

#### Classification

Finally, we learn machine learning models for each malware family to classify malwares. The scores and number of alignments calculated as mentioned above are used as the features in our classifiers.

We experimented with a number of machine learning models including logistic regression, support vector machines (SVM), decision tree and deep learning. Since the results did not vary significantly across models (see III in Results), we use logistic regression as our primary model because of its simplicity and interpretability.

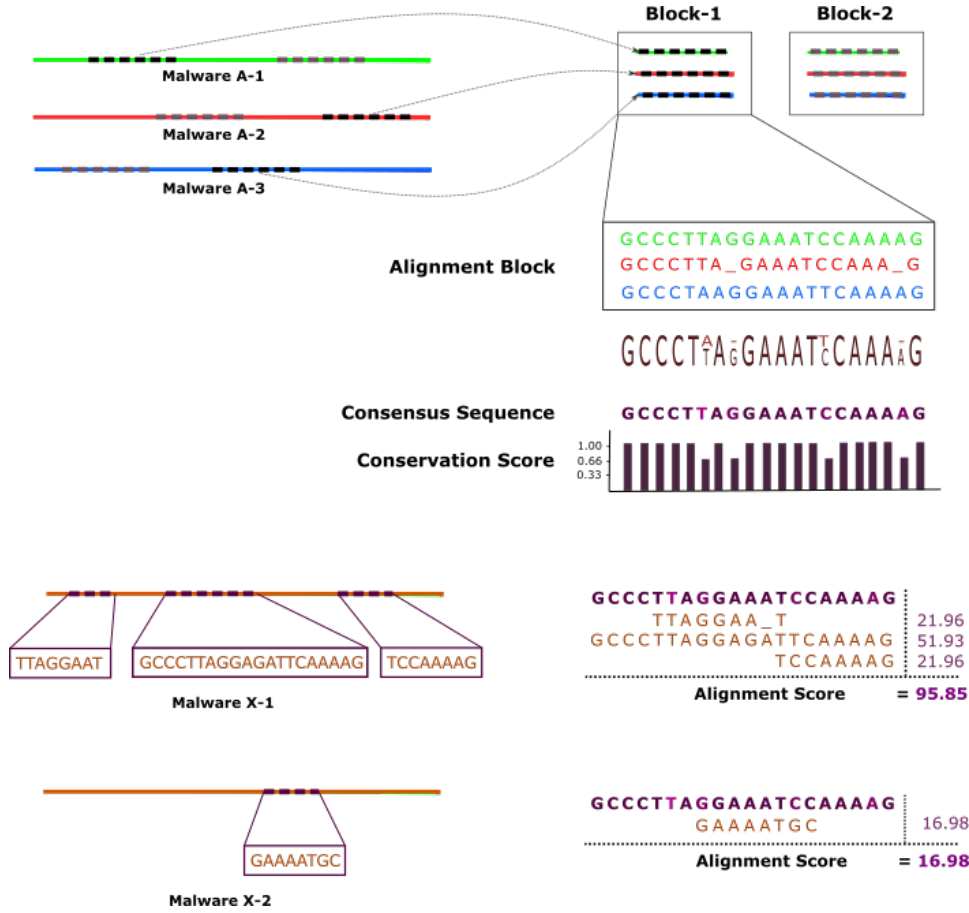


Fig. 2. Details of consensus sequence, conservation score, and alignment score generation from alignment blocks.

After the training phase, we get classifiers that can be used to classify or detect new instances as shown in Figure 1. The scores of the train and the test examples are calculated in the same way. The scores for new instances are passed to the classifiers to classify them into positive and negative instances. If a new sample is classified as negative by classifiers for all families, it can be considered as a benign sample.

## V. RESULTS

In the following sections, we first discuss the datasets used in this paper and subsequently present the results on these datasets.

### Datasets

*The Kaggle Microsoft Malware Classification Challenge (Big 2015):* The Kaggle Microsoft Malware Classification Challenge (Big 2015) [49] aimed to organize polymorphic malwares into 9 separate classes of malicious programs at a high level (see Table II). This challenge simulates the file input data processed on over 160 million computers by Microsoft’s real-time anti-malware detection products inspecting over 700 million computers per month.

Microsoft provided almost half a terabyte of input training and classification input data when uncompressed. They included:

Family Name	No of Train Samples	Type
Ramnit	1541	Worm
Lollipop	2478	Adware
Kelihos_ver3	2942	Backdoor
Vundo	475	Trojan
Simda	42	Backdoor
Tracur	751	TrojanDownloader
Kelihos_ver1	398	Backdoor
Obfuscator.ACY	1228	Any kind of obfuscated malware
Gatak	1013	Backdoor

TABLE II  
MALWARE FAMILIES IN THE KAGGLE DATASET

- 1) *Binary Files:* 10,868 training files containing the raw hexadecimal representation of the file’s binary content.
- 2) *Assembly Files:* 10,868 training files containing data extracted by the Interactive Disassembler (IDA) Tool. This information includes assembly command sequences, function calls and more.
- 3) *Training Labels:* Each training file name is a MD5 hash of the actual program. Each MD5 hash and the malware class it maps to are stored in the training label file.

From this, we constructed 9 balanced datasets for binary classification consisting of equal number of positive and

negative samples for each of the 9 malware families. In each dataset, the positive examples are all the samples from the corresponding family and the negative examples were chosen by randomly sampling from the 8 other families in the dataset. Then 20% of each dataset was set aside as the test sets while the remaining 80% was used as the training sets. For the machine learning approaches that require hyper-parameter selection, the 80% was further split into training (60%) and validation sets (20%).

*Microsoft Machine Learning Security Evasion Competition (2020) (MLSec) Dataset:* While the Kaggle Microsoft Malware Classification Challenge (Big 2015) dataset is a large and widely studied one, often malwares families only have a few samples - especially when they emerge initially. Therefore it is important to assess the performance of the methods on datasets with small number of instances per family. So, we applied our method on the Microsoft Machine Learning Security Evasion Competition (2020) [51] (MLSec) dataset, too. Here we used the dataset from ‘Defender Challenge’ which consisted of malware bytes code and their variants. Defenders’ challenge was to create a solution model that can defend against evasive variants created by the attackers. 49 malwares along with their evasive variants (submitted by the attackers) were found in this ‘Defender Challenge’ dataset. This dataset contains 49 original malwares with unique id from ‘01’ to ‘49’. Each malware contains a different number of evasive variants varying from 5 to 20. On average, a malware has 12 variants in this dataset. Similarly to the Kaggle Microsoft Malware dataset, 49 datasets were created which were then split into training, validation and test sets. During the split of this dataset, we always kept the original sample in the train set and the variants in the test set for each family, so that the test set can be considered as an evolution of the train set.

#### *Evaluation of machine learning algorithms*

First we assess the performance of various machine learning approaches on the Kaggle Microsoft Malware (Big 2015) dataset. The binary files of this dataset were converted to nucleotide sequence files and labelled using ‘Training Labels’ data. Then we generated common alignment blocks using SibeliaZ and constructed the consensus sequences as discussed in Methods. We generated the conservation scores for each consensus sequences using frequency of nucleotides which were then used as features in the machine learning models.

We experimented with logistic regression, decision trees and support vector machines (SVM). Table III shows the train and test accuracy for 80%-20% train-test split on Kaggle Microsoft Malware Classification Challenge (Big 2015) Dataset. We were unable to align instances of the ‘Lollipop’ family by SibeliaZ possibly due to the limitation of computational resources. Hence, the family was removed from our analysis. We observe that the algorithms show similar performances in terms of accuracy. So, we selected logistic regression for future experiments because of its simplicity and interpretability. We experimented with the hyper-parameters of logistic regression

and found that it gave the best results for ‘elasticnet’ penalty,  $C=0.05$  (regularization factor), ‘saga’ solver and  $l1\_ratio=0.5$ .

#### *Comparison with existing approaches*

Next we compare the performance of MALIGN with that of state of the art approaches, MalConv [33] (a neural network based approach using raw byte sequence), Ahmadi et al. [32] (a feature fusion based approach using byte and assembly files), and M-CNN [27] (a convolutional neural network (CNN) based approach relying on conversion to images) on the Kaggle Microsoft Malware (Big 2015) dataset. It is worth noting that models with multiclass loss as low as 0.00283 have been reported for this dataset. However, we compare with MalConv and M-CNN, as they have been successfully applied to many different datasets. We compared MALIGN with Ahmadi et al.’s Feature-Fusion method, because to our knowledge, this was the closest to the accuracy of the winning team of the Kaggle competition.

We also implement a deep learning based approach that classifies malwares using the alignment scores calculated by MALIGN. The architectures of the deep learning based approach on alignment scores as well as architectures of MalConv and M-CNN are shown in Figure 3.

The training and test accuracy of MALIGN with logistic regression and deep learning along with those of Feature-Fusion, MalConv and M-CNN are shown in Tables IV and V. Table V shows that MALIGN (Deep Learning) has better accuracy on the test set than other approaches. MALIGN (Deep Learning) has the best accuracy of 98.59%. Table IV shows that, on train set, MALIGN has better accuracy than MalConv and M-CNN, and the difference with Feature-Fusion is negligible.

#### *Running Time*

We run the experiments of different methods on different platforms. However, to provide an idea, the total running time (from data processing to classification) for MLSEC dataset is given in the table VI. Although inference time grows linearly with the number of families, once the signatures are found in MALIGN, they can be trimmed if needed (based on how conserved they are) and then the alignment of the sequence of the new instance and the signature can be sped up.

Ahmadi et al. [32] Feature-Fusion method has not been included in the table VI because running it on Mlsec dataset was not possible due to its limitation of using both the byte and the assembly file. But we can get an estimation of its running time on the Kaggle Microsoft dataset from the paper [32]. For example, it takes almost 17 hours 15 minutes to extract ‘REG’ feature from all samples, and it extracts in total of 14 features.

#### *Applicability with limited amount of data and features*

Although deep learning based approaches have been widely applied for malware classification and detection, they require extensive amount of data for training and tend to overfit in absence of that. Tables IV and V show that Feature-Fusion, MalConv and M-CNN perform well for most malware



Family Name	Train Accuracy			Test Accuracy		
	Logistic regression	Decision tree	SVM	Logistic regression	Decision tree	SVM
Ramnit	99.91	99.91	99.91	99.64	99.82	99.64
Kelihos_ver3	99.83	99.94	99.81	99.27	99.27	99.27
Vundo	100	100	100	97.4	97.4	97.4
Simda	100	100	97.92	84.62	84.62	84.62
Tracur	100	100	99.5	97.3	94.6	96.3
Kelihos_ver1	100	100	99.5	96.7	98.9	98.9
Obfuscator.ACY	100	100	100	95.9	92.7	96
Gatak	99.17	99.17	99.17	96.2	96.2	96.2
Overall	99.82	99.86	99.74	97.99	97.42	98.02

TABLE III

CLASSIFICATION ACCURACY ON KAGGLE MICROSOFT MALWARE (BIG 2015) DATASET FOR DIFFERENT MACHINE LEARNING MODELS

Family Name	MALIGN (Logistic Regression)	MALIGN (Deep Learning)	Feature-Fusion	MalConv	M-CNN
Ramnit	99.91	99.58	100	98.39	97.64
Kelihos_ver3	99.83	99.86	100	99.89	99.91
Vundo	100	98.26	100	99.4	99.26
Simda	100	94.44	100	100	100
Tracur	100	98.18	100	98.74	98.43
Kelihos_ver1	100	98.92	100	98.54	99.52
Obfuscator.ACY	100	98.66	100	97.33	99.70
Gatak	99.17	99.2	100	99.15	92.69
Overall	99.82	99.24	100	98.96	98.4

TABLE IV

PERFORMANCE OF DIFFERENT MODELS ON TRAIN-SET OF KAGGLE MICROSOFT MALWARE CLASSIFICATION CHALLENGE DATASET (BIG 2015)

Family Name	MALIGN (Logistic Regression)	MALIGN (Deep Learning)	Feature-Fusion	MalConv	M-CNN
Ramnit	99.64	99.46	98.7	95.66	88.44
Kelihos_ver3	99.27	99.82	99.18	100	99.72
Vundo	97.4	98.70	95.79	94.89	97.21
Simda	84.62	84.62	76.47	52.94	62.5
Tracur	97.3	98.2	98.34	93.91	94.55
Kelihos_ver1	96.7	95.7	98.75	96.08	94.67
Obfuscator.ACY	95.9	96.39	98.98	94.42	91.74
Gatak	96.2	98.37	98.03	98.67	88.68
Overall	97.99	98.59	98.52	96.95	94.10

TABLE V

PERFORMANCE OF DIFFERENT MODELS ON TEST-SET OF KAGGLE MICROSOFT MALWARE CLASSIFICATION CHALLENGE DATASET (BIG 2015)

Method	Running Time
MALIGN	18hr 4min
MalConv	18hr 13min
M-CNN	14hr 26min

TABLE VI

MODEL RUNNING TIME FOR DIFFERENT METHODS ON MLSEC DATASET

families. However, we observe that, for Type 5 (Simda), which has only 42 samples, the test accuracy of Feature-Fusion, MalConv and M-CNN are only 76.47%, 52.94% and 62.5% respectively whereas MALIGN has 84.62% test accuracy. The other methods having training accuracy of 100% indicates overfitting. Same observation can be made for Type 4 (Vundo) that has second smallest number of train samples.

Moreover, MALIGN needs only the raw byte sequence whereas the Feature-Fusion method needs byte sequence, assembly file, address for byte sequence, section information from PE etc. So, even when only the binary executable files are available MALIGN will work perfectly, but typical approaches like Feature-Fusion [32] that needs various features, will have

to use other tools (such as IDA) to work properly.

We also compare performances of the methods on the MLSec dataset (Microsoft Machine Learning Security Evasion Competition (2020) [51]). Since this dataset contains limited number of variants created in almost real-time, this can be used to identify how our method works on zero-day malwares when only a limited number of samples are available.

Because of the limited number of instances in this dataset, the validation set is very small for some types. So we run the deep learning models on a 80%-20% train-test split as well as 60%-20%-20% train-validation-test split of the data. Tables VII and VIII summarizes the performances of the models on all 49 types for train-test and train-validation-test split respectively whereas Figures 4a and 4b provide radar charts showing training and test accuracy in all 49 types individually.



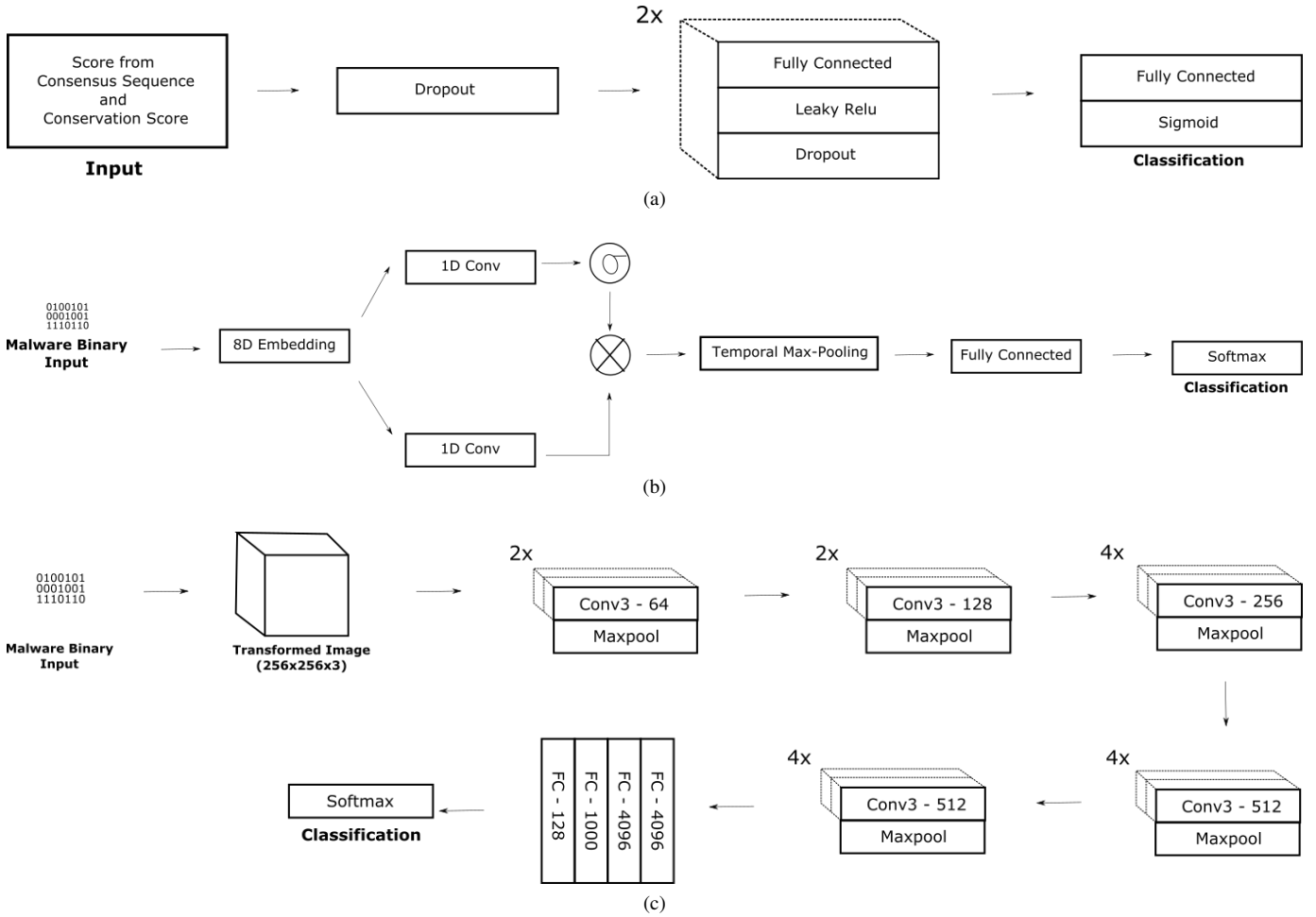


Fig. 3. Architectures of (a) deep learning model on alignment scores, (b) the MalConv model [33], and (c) the M-CNN model [27].

	Train Accuracy	Test Accuracy
MALIGN(Logistic Regression)	98.18	<b>80.42</b>
MALIGN(Deep Learning)	97.72	80.00
MalConv	<b>98.24</b>	79.22
M-CNN	96.99	71.09

TABLE VII

PERFORMANCE OF MODELS ON MLSEC DATASET FOR TRAIN-TEST SPLIT

	Train Accuracy	Validation Accuracy	Test Accuracy
MALIGN (Deep Learning)	97.53	<b>81.67</b>	<b>79.58</b>
MalConv	95.39	81.22	64.45
M-CNN	<b>98.4</b>	74.29	73.83

TABLE VIII

PERFORMANCE OF MODELS ON MLSEC DATASET FOR TRAIN-VALIDATION-TEST SPLIT

We observe that MALIGN outperforms the other deep learning based models overall on the MLSec dataset regardless of the splitting. This highlights the advantage of explicit identification of critical code blocks when data is limited.

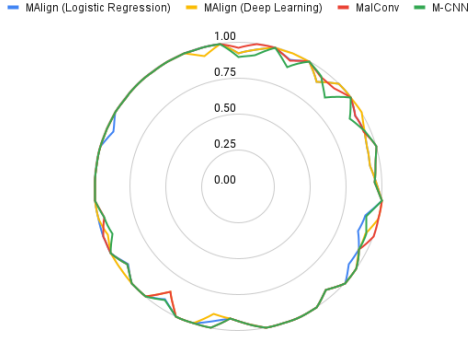
From Figures 4a and 4b, we observe that on the train set, all approaches have consistent performance, but on the test set M-CNN and MalConv is relatively inconsistent. For example,

type 29 and 43 have 10 and 5 available variants respectively, and M-CNN's test accuracy is 0 on both.

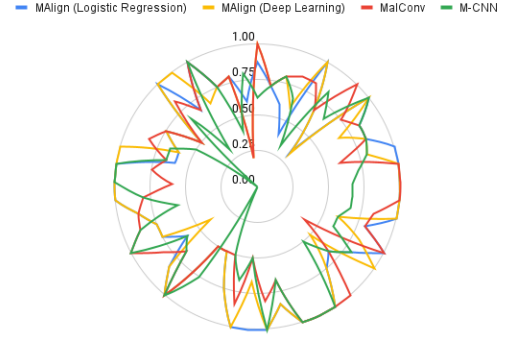
#### Robustness to adversarial attacks

A major issue with deep learning based malware detection approaches is - they can often be evaded by adding selective content to it. In deep learning based methods, the gradient attack can be used to find such selected content. Since MALIGN relies on finding conserved blocks critical to the malware families for classification through sequence alignment, score calculation, and an interpretable logistic regression model, it should in its principle to be inherently robust to such attacks.

We tried to investigate the robustness of our method compared to conventional malware detection techniques. We used the evasion technique on MalConv model that was proposed by Kolosnjaji et al. [1]. It creates adversarial samples just by modifying (or padding) approximately 1.25% of the total size of malwares which can successfully evade the MalConv model with high percentage. The evasion rate increases with the percentage of modification on malware samples. We experimented with their implementation [52] on some of the types in the MLSec dataset, and the results are shown in Table IX.



(a) Radar chart showing accuracy for MLSec-Train Dataset



(b) Radar chart showing accuracy for MLSec-Test Dataset

Fig. 4. Radar chart showing accuracy of different models on each of the 49 malware types along the perimeter on MLSec-Train Dataset

Type	MalConv			MALIGN
	Evaded/Total Train set	Test set	Evasion Rate	Evasion Rate
1	5/15	2/4	36.84%	0.00%
11	1/9	1/2	18.18%	0.00%
12	5/13	2/4	41.18%	0.00%
28	5/8	2/2	70%	0.00%
45	4/5	2/2	85.71%	0.00%

TABLE IX

GRADIENT ATTACK RESULTS ON MALCONV ON SOME TYPES IN THE MLSEC DATASET

We then applied MALIGN on these evasive samples and all of them were successfully detected with almost 100% prediction probability. It is worth mentioning that since the evasion technique did not have access to the MALIGN model, the experiment is not rigorous. However, since MALIGN finds the blocks critical for function in a malware family and classifies them based on those, to evade MALIGN, the malware attackers will have to go through an incommodious process of changing those blocks without changing its intended features and semantics. Moreover, gradient attack based evasion techniques cannot directly be applied on MALIGN because we are not feeding the malware file directly to our model unlike many other techniques. Besides gradient attack, other typical obfuscation techniques will not be able to evade MALIGN. Some of those have been discussed below-

- *Including pieces of other malware to confuse family detection:* The inclusion of pieces from other malware will be detected but code-pieces of the original family will also be detected at the same time, and consequently, the score for the original family will be higher since it contains way more code-pieces than the other family. Thus, MALIGN will still identify that sample correctly.
- *Intersperse instructions into other benign programs to dilute the signal:* Interspersed instructions will still be detected as sequence alignment works in the presence of substitutions, insertions, and deletions of nucleotides (instructions). Now the prediction of MALIGN will depend on the score and number of occurrences of these interspersed instructions.

- *Using indirect addressing:* Use of indirect addressing may cause some mismatches, but still there will be some matching because the source and destination register will have to be the same. Moreover, there will still be other preserved instructions except addresses to protect the malware semantics which will be captured by MALIGN.

There are plausible approaches to attack the present implementation which can be addressed using classifiers with only non-negative weights and other techniques in the future.

## VI. CASE STUDIES

An important advantage of MALIGN is - it is interpretable and insights can be derived about malware families through a simple *backtracking* process. The process is as follows:

- Find the blocks that are assigned high weights by the logistic regression model
- Select the blocks that are highly conserved from the above list
- Process the alignment file (MAF) to determine the sequences and their indices that constructed the blocks
- Locate the code fragments corresponding to the sequences found in Step (iii)

This process can be used to uncover potential malicious code as discussed next.

### Detection of obfuscated malicious code

Different techniques and methods are used by attackers or malware creators to obfuscate malicious, harmful code segments to evade anti-malware tools. MALIGN, our proposed method gives us the ability to find responsible malicious code segments from malware files by the backtracking process sketched above and analyze the blocks of code that are highly conserved.

We backtracked from the aligned blocks to the assembly code on some randomly selected samples from the Kaggle Microsoft Malware (Big 2015) dataset. In some cases, we found evidence of malware obfuscation. Figure 5 is an example of such case (data-transformation obfuscation technique in this case). Figure 5a, 5b and 5c are snapshots of three different

gGPBDhr24EI875LfTNKV 10001230 10001250

```
.text:10001233
.text:10001233
.text:10001233
.text:10001233
.text:10001233 33 C0
.text:10001235 EB 17
.text:10001237
.text:10001237 9C
.text:10001238 A5
.text:10001239
.text:10001239
.text:10001239 EB 2B
.text:10001239
.text:1000123B 7A
.text:1000123C 2B 88 21 46 07 34 5D D2 A3 A0 59 1E FF CC 15 2A
.text:1000124C 1B B8
.text:1000124E
.text:1000124E
.text:1000124E
.text:1000124E EB E9
.text:1000124E
.text:10001250 91 F6 F7 64
```

```
; ===== SUBROUTINE =====
sub_10001233 proc far ; CODE XREF: DllEntryPoint+31FD84p
xor eax, eax
jmp short loc_1000124E
; -----
pushf
movsd
loc_10001239: ; CODE XREF: sub_10001233:loc_1000124E84j
jmp short loc_10001266
; -----
db 7Ah
dd 4621882Bh, 0D25D3407h, 1E59A0A3h, 2A15CCFFh
db 1Bh, 0B8h
; -----
loc_1000124E: ; CODE XREF: sub_10001233+2C4Bj
jmp short loc_10001239
; -----
dd 64F7F691h
```

(a) Code in .text segment

```
2mYSX6J9Dd51kpLuNmC8 10004200 10004220
```

```
.rdata:10004204 EB
.rdata:10004205 17
.rdata:10004206 9C
.rdata:10004207 A5
.rdata:10004208 EB
.rdata:10004209 2B
.rdata:1000420A 7A
.rdata:1000420B 2B
.rdata:1000420C 88
.rdata:1000420D 21
.rdata:1000420E 46
.rdata:1000420F 07
.rdata:10004210 34
.rdata:10004211 5D
.rdata:10004212 D2
.rdata:10004213 A3
.rdata:10004214 A0
.rdata:10004215 59
.rdata:10004216 1E
.rdata:10004217 FF
.rdata:10004218 CC
.rdata:10004219 15
.rdata:1000421A 2A
.rdata:1000421B 1B
.rdata:1000421C B8
.rdata:1000421D EB
.rdata:1000421E E9
.rdata:1000421F 91
.rdata:10004220 F6
.rdata:10004221 F7
.rdata:10004222 64
```

```
db 0EBh ; e
db 17h
db 9Ch ; e
db 0A5h ; Y
db 0EBh ; e
db 2Bh ; +
db 7Ah ; z
db 2Bh ; +
db 88h ; ^
db 21h ; !
db 46h ; F
db 7
db 34h ; 4
db 5Dh ; j
db 0D2h ; 0
db 0A3h ; f
db 0A0h ;
db 59h ; Y
db 1Eh
db 0FFh
db 0CCh ; i
db 15h
db 2Ah ; *
db 1Bh
db 0B8h ;
db 0EBh ; e
db 0E9h ; e
db 91h ; `
db 0F6h ; o
db 0F7h ; +
db 64h ; d
```

```
kRUx3TuoJSgp0sqDzNGX 10005E90 10005EB5
```

```
.data:10005E90 3B
.data:10005E91 C0
.data:10005E92 23
.data:10005E93 A7
.data:10005E94 87
.data:10005E95 3C
.data:10005E96 24
.data:10005E97 EB
.data:10005E98 17
.data:10005E99 9C
.data:10005E9A A5
.data:10005E9B EB
.data:10005E9C 2B
.data:10005E9D 7A
.data:10005E9E 2B
.data:10005E9F 88
.data:10005EA0 21
.data:10005EA1 46
.data:10005EA2 07
.data:10005EA3 34
.data:10005EA4 5D
.data:10005EA5 D2
.data:10005EA6 A3
.data:10005EA7 A0
.data:10005EA8 59
.data:10005EA9 1E
.data:10005EAA FF
.data:10005EAB CC
.data:10005EAC 15
.data:10005EAD 2A
.data:10005EAE 1B
.data:10005EAF B8
.data:10005EB0 EB
.data:10005EB1 E9
.data:10005EB2 91
.data:10005EB3 F6
.data:10005EB4 F7
.data:10005EB5 64
```

```
db 3Bh ; ;
db 0C0h ; A
db 23h ; #
db 0A7h ; S
db 87h ; +
db 3Ch ; <
db 24h ; $
db 0EBh ; e
db 17h
db 9Ch ; e
db 0A5h ; Y
db 0EBh ; e
db 2Bh ; +
db 7Ah ; z
db 2Bh ; +
db 88h ; ^
db 21h ; !
db 46h ; F
db 7
db 34h ; 4
db 5Dh ; j
db 0D2h ; 0
db 0A3h ; f
db 0A0h ;
db 59h ; Y
db 1Eh
db 0FFh
db 0CCh ; i
db 15h
db 2Ah ; *
db 1Bh
db 0B8h ;
db 0EBh ; e
db 0E9h ; e
db 91h ; `
db 0F6h ; o
db 0F7h ; +
db 64h ; d
```

(b) Code in .rdata segment

(c) Code in .data segment

Fig. 5. Code obfuscation in data segment (a) A code fragment in a malware sample, (b) and (c) Same code obfuscated in .rdata and .data segments indicated by the hex-codes

malware files from ‘Vundo’ malware family of ‘Trojan’ type. All samples contain the same hex-code but in different segments. Figure 5b and 5c conceal the same code string of Figure 5a in its .rdata and .data section with db respectively. In short, the same machine code was transformed to its hex form

and placed into data section, possibly to evade anti-malware tools.

## VII. CONCLUSIONS

In this paper, we presented a malware detection tool MALIGN based on a recently developed multiple whole-genome alignment tool, SibeliaZ [20]. Sequence alignment based approaches have been used for malware analysis in the past, but the use of a whole-genome alignment tool makes MALIGN scalable to long malware sequences and protects against trivial adversarial attacks such as code obfuscation. The method is also interpretable and can be used to derive insights on malwares such as identification of critical code blocks and code obfuscation.

We have applied MALIGN on the Kaggle Microsoft Malware Classification Challenge (Big 2015) and the Microsoft Machine Learning Security Evasion Competition (2020) (MLSec) datasets, and observed that it outperforms widely used deep learning based methods such as MalConv, M-CNN, and feature fusion method (Ahmadi et al.).

Preliminary experiments show that MALIGN is robust to common adversarial attacks such as padding and modification of sequences. In future, the method may be tested against other possible types evasion techniques and the machine learning algorithms can be adjusted accordingly. In addition, other whole-genome alignment tools such as Progressive Cactus [19] may be experimented with.

## REFERENCES

- [1] B. Kolosnjaji, A. Demontis, B. Biggio, D. Maiorca, G. Giacinto, C. Eckert, and F. Roli, "Adversarial malware binaries: Evading deep learning for malware detection in executables," in *2018 26th European signal processing conference (EUSIPCO)*. IEEE, 2018, pp. 533–537.
- [2] K. Grosse, N. Papernot, P. Manoharan, M. Backes, and P. McDaniel, "Adversarial perturbations against deep neural networks for malware classification," *arXiv preprint arXiv:1606.04435*, 2016.
- [3] A. Huang, A. Al-Dujaili, E. Hemberg, and U.-M. O'Reilly, "On visual hallmarks of robustness to adversarial malware," *arXiv preprint arXiv:1805.03553*, 2018.
- [4] F. Kreuk, A. Barak, S. Aviv-Reuven, M. Baruch, B. Pinkas, and J. Keshet, "Deceiving end-to-end deep learning malware detectors using adversarial examples," *arXiv preprint arXiv:1802.04528*, 2018.
- [5] H. S. Anderson, A. Kharkar, B. Filar, D. Evans, and P. Roth, "Learning to evade static pe machine learning malware models via reinforcement learning," *arXiv preprint arXiv:1801.08917*, 2018.
- [6] A. Al-Dujaili, A. Huang, E. Hemberg, and U.-M. O'Reilly, "Adversarial deep learning for robust detection of binary encoded malware," in *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2018, pp. 76–82.
- [7] J. W. Stokes, D. Wang, M. Marinescu, M. Marino, and B. Bussone, "Attack and defense of dynamic analysis-based, adversarial neural malware detection models," in *MILCOM 2018-2018 IEEE Military Communications Conference (MILCOM)*. IEEE, 2018, pp. 1–8.
- [8] W. Hu and Y. Tan, "Black-box attacks against rnn based malware detection algorithms," in *Workshops at the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [9] I. Rosenberg, A. Shabtai, L. Rokach, and Y. Elovici, "Generic black-box end-to-end attack against state of the art api call based malware classifiers," in *International Symposium on Research in Attacks, Intrusions, and Defenses*. Springer, 2018, pp. 490–510.
- [10] W. Hu and Y. Tan, "Generating adversarial malware examples for black-box attacks based on gan," *arXiv preprint arXiv:1702.05983*, 2017.
- [11] S. B. Needleman and C. D. Wunsch, "A general method applicable to the search for similarities in the amino acid sequence of two proteins," *Journal of Molecular Biology*, vol. 48, no. 3, pp. 443–453, 1970.
- [12] T. F. Smith and M. S. Waterman, "Identification of common molecular subsequences," *Journal of Molecular Biology*, vol. 147, no. 1, pp. 195–197, 1981.
- [13] J. D. Thompson, D. G. Higgins, and T. J. Gibson, "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Research*, vol. 22, no. 22, pp. 4673–4680, 1994.
- [14] K. Katoh, K. Misawa, K.-i. Kuma, and T. Miyata, "MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform," *Nucleic Acids Research*, vol. 30, no. 14, pp. 3059–3066, 2002.
- [15] R. C. Edgar, "MUSCLE: multiple sequence alignment with high accuracy and high throughput," *Nucleic Acids Research*, vol. 32, no. 5, pp. 1792–1797, 2004.
- [16] A. L. Delcher, S. Kasif, R. D. Fleischmann, J. Peterson, O. White, and S. L. Salzberg, "Alignment of whole genomes," *Nucleic Acids Research*, vol. 27, no. 11, pp. 2369–2376, 1999.
- [17] A. C. Darling, B. Mau, F. R. Blattner, and N. T. Perna, "Mauve: multiple alignment of conserved genomic sequence with rearrangements," *Genome Research*, vol. 14, no. 7, pp. 1394–1403, 2004.
- [18] B. Paten, D. Earl, N. Nguyen, M. Diekhans, D. Zerbino, and D. Haussler, "Cactus: Algorithms for genome multiple sequence alignment," *Genome Research*, vol. 21, no. 9, pp. 1512–1528, 2011.
- [19] J. Armstrong, G. Hickey, M. Diekhans, I. T. Fiddes, A. M. Novak, A. Deran, Q. Fang, D. Xie, S. Feng, J. Stiller et al., "Progressive Cactus is a multiple-genome aligner for the thousand-genome era," *Nature*, vol. 587, no. 7833, pp. 246–251, 2020.
- [20] I. Minkin and P. Medvedev, "Scalable multiple whole-genome alignment and locally collinear block construction with SibeliaZ," *Nature communications*, vol. 11, no. 1, pp. 1–11, 2020.
- [21] C. Wressnegger, K. Freeman, F. Yamaguchi, and K. Rieck, "Automatically inferring malware signatures for anti-virus assisted attacks," in *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, 2017, pp. 587–598.
- [22] M. Zakeri, F. Faraji Daneshgar, and M. Abbaspour, "A static heuristic approach to detecting malware targets," *Security and Communication Networks*, vol. 8, no. 17, pp. 3015–3027, 2015.
- [23] L. Martignoni, E. Stinson, M. Fredrikson, S. Jha, and J. C. Mitchell, "A layered architecture for detecting malicious behaviors," in *International Workshop on Recent Advances in Intrusion Detection*. Springer, 2008, pp. 78–97.
- [24] C. Willems, T. Holz, and F. Freiling, "Toward automated dynamic malware analysis using cwsandbox," *IEEE Security & Privacy*, vol. 5, no. 2, pp. 32–39, 2007.
- [25] M. G. Schultz, E. Eskin, F. Zadok, and S. J. Stolfo, "Data mining methods for detection of new malicious executables," in *Proceedings 2001 IEEE Symposium on Security and Privacy. S&P 2001*. IEEE, 2000, pp. 38–49.
- [26] L. Nataraj, S. Karthikeyan, G. Jacob, and B. S. Manjunath, "Malware images: visualization and automatic classification," in *Proceedings of the 8th international symposium on visualization for cyber security*, 2011, pp. 1–7.
- [27] M. Kalash, M. Rochan, N. Mohammed, N. D. Bruce, Y. Wang, and F. Iqbal, "Malware classification with deep convolutional neural networks," in *2018 9th IFIP International Conference on New Technologies, Mobility and Security (NTMS)*. IEEE, 2018, pp. 1–5.
- [28] R. K. Shahzad, N. Lavesson, and H. Johnson, "Accurate adware detection using opcode sequence extraction," in *2011 Sixth International Conference on Availability, Reliability and Security*. IEEE, 2011, pp. 189–195.
- [29] R. Lu, "Malware detection with lstm using opcode language," *arXiv preprint arXiv:1906.04593*, 2019.
- [30] I. Santos, J. Devesa, F. Brezo, J. Nieves, and P. G. Bringas, "Opem: A static-dynamic approach for machine-learning-based malware detection," in *International joint conference CISIS'12-ICEUTE' 12-SOCO' 12 special sessions*. Springer, 2013, pp. 271–280.
- [31] J. Yan, Y. Qi, and Q. Rao, "Detecting malware with an ensemble method based on deep neural network," *Security and Communication Networks*, vol. 2018, 2018.
- [32] M. Ahmadi, D. Ulyanov, S. Semenov, M. Trofimov, and G. Giacinto, "Novel feature extraction, selection and fusion for effective malware family classification," in *Proceedings of the sixth ACM conference on data and application security and privacy*, 2016, pp. 183–194.
- [33] E. Raff, J. Barker, J. Sylvester, R. Brandon, B. Catanzaro, and C. K. Nicholas, "Malware detection by eating a whole exe," in *Workshops at*

the Thirty-Second AAAI Conference on Artificial Intelligence, 2018, pp. 268–276.

- [34] T. Bai, J. Luo, J. Zhao, B. Wen, and Q. Wang, “Recent advances in adversarial training for adversarial robustness,” *arXiv preprint arXiv:2102.01356*, 2021.
- [35] K. Grosse, N. Papernot, P. Manoharan, M. Backes, and P. McDaniel, “Adversarial examples for malware detection,” in *European symposium on research in computer security*. Springer, 2017, pp. 62–79.
- [36] Y. Zhang, H. Li, Y. Zheng, S. Yao, and J. Jiang, “Enhanced dnns for malware classification with gan-based adversarial training,” *Journal of Computer Virology and Hacking Techniques*, vol. 17, no. 2, pp. 153–163, 2021.
- [37] A. Al-Dujaili, A. Huang, E. Hemberg, and U.-M. O’Reilly, “Adversarial deep learning for robust detection of binary encoded malware,” in *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2018, pp. 76–82.
- [38] J. Cohen, E. Rosenfeld, and Z. Kolter, “Certified adversarial robustness via randomized smoothing,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 1310–1320.
- [39] Y. Chen, S. Wang, D. She, and S. Jana, “On training robust {PDF} malware classifiers,” in *29th {USENIX} Security Symposium ({USENIX} Security 20)*, 2020, pp. 2343–2360.
- [40] Í. Íncar Romeo, M. Theodorides, S. Afroz, and D. Wagner, “Adversarially robust malware detection using monotonic classification,” in *Proceedings of the Fourth ACM International Workshop on Security and Privacy Analytics*, 2018, pp. 54–63.
- [41] Y. Chen, A. Narayanan, S. Pang, and B. Tao, “Multiple sequence alignment and artificial neural networks for malicious software detection,” in *2012 8th International Conference on Natural Computation*. IEEE, 2012, pp. 261–265.
- [42] A. Narayanan, Y. Chen, S. Pang, and B. Tao, “The effects of different representations on malware motif identification,” in *2012 Eighth International Conference on Computational Intelligence and Security*. IEEE, 2012, pp. 86–90.
- [43] V. Naidu and A. Narayanan, “Further experiments in biocomputational structural analysis of malware,” in *2014 10th International Conference on Natural Computation (ICNC)*. IEEE, 2014, pp. 605–610.
- [44] D. Kirat and G. Vigna, “Malgene: Automatic extraction of malware analysis evasion signature,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015, pp. 769–780.
- [45] V. Naidu and A. Narayanan, “Needleman-wunsch and smith-waterman algorithms for identifying viral polymorphic malware variants,” in *2016 IEEE 14th Intl Conf on Dependable, Autonomic and Secure Computing, 14th Intl Conf on Pervasive Intelligence and Computing, 2nd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech)*. IEEE, 2016, pp. 326–333.
- [46] I. K. Cho, T. G. Kim, Y. J. Shim, M. Ryu, and E. G. Im, “Malware analysis and classification using sequence alignments,” *Intelligent Automation & Soft Computing*, vol. 22, no. 3, pp. 371–377, 2016.
- [47] H. Kim, J. Kim, Y. Kim, I. Kim, K. J. Kim, and H. Kim, “Improvement of malware detection and classification using api call sequence alignment and visualization,” *Cluster Computing*, vol. 22, no. 1, pp. 921–929, 2019.
- [48] J. Drew, T. Moore, and M. Hahsler, “Polymorphic malware detection using sequence classification methods,” in *2016 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2016, pp. 81–87.
- [49] “Microsoft malware classification challenge (big 2015),” accessed July 7, 2019 =<https://www.kaggle.com/c/malwareclassification/>.
- [50] J. D. Kececioglu and E. W. Myers, “Combinatorial algorithms for dna sequence assembly,” *Algorithmica*, vol. 13, no. 1, pp. 7–51, 1995.
- [51] “Microsoft mmachine learning security evasion competition (2020),” <https://mlsec.io/>.
- [52] <https://github.com/yuxiaorun/MalConv-Adversarial>.