

BERT - Based Detection and Classification of Algorithmically Generated Domains in Malware Communication

Santosh Kumar Kande

USA

ABSTRACT

Malware threats. Particularly ransomware, often employ Algorithmically Generated Domains (AGD) for communication with Command & Control (C&C) servers. This paper introduces a novel approach for AGD detection using the Bidirectional Encoder Representations from Transformer (BERT) model, eliminating the need for intricate feature selection or hyperparameter tuning. The proposed method effectively addresses the challenge posed by sophisticated domain generation techniques, including dictionary-based and random character approaches. Experimental results demonstrate the superior performance of the BERT model in both AGD detection and classification tasks, achieving a precision, recall, and accuracy score of 0.99. The approach proves effective against diverse domain generation algorithms, enhancing the current state-of-the-art methods for securing networks against evolving malware threats.

*Corresponding author

Santosh Kumar Kande, USA.

Received: March 11, 2024; **Accepted:** March 14, 2024; **Published:** March 25, 2024

Keywords: Domain Generation Algorithm, Domain Name Server, Malware, NLP, Transformer

Introduction

Malware, such as ransomware, must interact with their respective Command & Control (C&C) server once infected. Malware relies on domain names produced by the Domain Generation Algorithm (DGA) to communicate with command and control servers. Attackers register Algorithmically Generated Domain (AGD) names. Attackers register Algorithmically Generated Domain (AGD) names. Upon successful contact, the malware either leaks data to the C&C server or infects additional machines in the victim's network. In a ransomware attack, malware encrypts the victim's data.

Malicious domain names, such as ocufxskoieqqv.com, might be gibberish. A huge number of lookup messages for such nonsensical domains, on the other hand, make it harder to discover and block rogue domains. DGAs now create domain names based on valid dictionary or wordlist names, such as scoreadmireluckapplyfitcouple.com, rather than random character-based names. Domain blacklisting is a common strategy for preventing such domain names from being contacted. However, this method can be easily evaded using sophisticated big wordlist-based domain names.

This paper presents the implementation of an AGD detection technique based on the Bidirectional Encoder Representational from Transformer (BERT) model. This approach doesn't require feature selection or hyperparameter tuning of model.

Related Work

Kührer, et al. found that public and vendor-provided domain name blocklists contained only 20% of major malware families and failed to provide any protection AGD names [1]. Antonakakis, et al. suggested clustering and classification approach to process the domain names in respective DGA families [2]. Woodbridge et al. used Long Short Term Memory (LSTM) architecture to detect the AGD names; however, faced challenges with dictionary or wordlist-based DGA [3]. Yu, et al. proposed a system for inline AGD detection at the DNS server using a Convolutional Neural Network (CNN) and LSTM-based approach [4]. A hybrid neural network approach comprising of parallel CNN and attention-based Bi-LSTM layers was designed by Yang, et al. For multiclass classification, they obtained a macro averaging F1 score of 0.7653 with proposed hybrid architecture [5]. Another similar approach utilizing attention based mechanism was proposed by Fangali, et al. that uses CNN, Bi-LSTM and Attention layer after the embedding layer to do the detection as well as multiclass classification of various DGAs [6]. They obtained micro average F1 score of 0.89 and macro average score of 0.83 on more than 20 real-world AGD name list and legitimate domain name lists, including word list based DGAs such as matsnu and supbbox.

Methodology

Starting with basic neural networks, increasingly sophisticated models such as CNN, RNN, LSTM, Bi-LSTM, Stacked-LSTM, and eventually BERT were developed Character level embedding was utilized in all of the studies.

DNN with 3 dense hidden layers and 128 neurons was utilized, with the relu activation function. Dropout was employed after each dense layer to avoid overfitting. The CNN model was made

up of a SpatialDropout1D layer, one convolution 1D layer, and a mac-pooling layer. Stacked LSTM units compromise two bi-directional LSTM layers, each with 64 units.

Transformers and Bidirectional Encoder Representations from Transformers (BERT)

Similar to the human brain, the attention mechanism in artificial neural networks aids in focusing on more relevant words in a phrase while determining how much significance to assign to each word. Vaswani et al. presented an attention-based Transformer neural network design. They eliminated the recurrent layers, which reduced training time and allowed them to use parallel processes [7]. Encoder and decoder components make up the transformer architecture. For neural machine translation, source language statement is embedded and processed by encoder, while the target language translation so far is produced by decoder as one token at a time. Attention mechanism works by querying a query vector against database of key vectors that are mapped to set of values (output values so far) as depicted in equation 1.

$$\text{Attention}(Q,K,V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

Attention is a scaled dot product that assigns a probability distribution to keys that are closely matching with the query vector Q. In practice Transformers employ multi-head attention, which entails numerous projections utilizing equation 1.

Devlin, et al. in 2018 used transfer learning to create Bidirectional Encoder Representations from Transformer (BERT) [8]. Transfer learning enables a model trained on one dataset to be used for a specific downstream task involving a different dataset. It is based on the Transformer architecture, although it only employs the Transformer encoder. BERT training is divided into two stages: Pre-training and Fine-Tuning.

During the **Pre-Training Phase**, BERT builds natural language understanding through unsupervised learning utilizing Masked Language Modeling (MLM) and Next Sentence Prediction (NSP). MLM is an unsupervised task in which a random word in a sentence is masked and the model predicts the probability distribution over output words. BERT employs a deep bidirectional model for language modeling rather than just left-to-right or concatenating two model outputs. In the case of NSP, given a pair of sentences A and B, if B follows A, it is labelled as *IsNext*; otherwise, it is labelled as *NotNext*. BERT is trained on the massive Wikipedia and BookCorpus.

During the **Fine-Tuning Phase**, BERT is trained on task-specific labelled datasets such as sentiment classification or question-answering datasets. During this step, the pre-trained model is trained from start to finish again, and all parameters are fine-tuned for the specific task. This is a pretty low-cost and less time-consuming activity.

For the problem at the hand, $BERT_{base}$ model with total of 110M parameters was used. Fine-tuning was done using labelled DGA dataset as shown in Figure 1.

Experimental Results and Discussion

Dataset

For the experimental investigation, real-world domain names generated by domain generation algorithms were utilized [9]. The dataset contained 110,000 genuine domain names as well as about 10000 domain names created by various DGAs. The dataset was divided into two parts: training and testing, with a 70:30 mix of randomly selected samples. Character level encoding was used to encode the entire dataset. The task was divided into two parts: DGA detection and DGA classification.

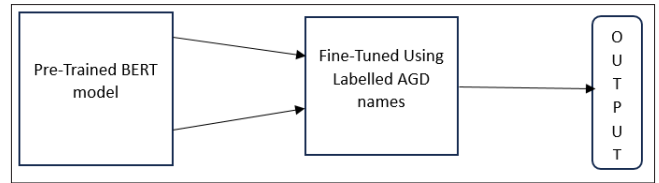


Figure 1: Fine Tuning of Pre-Trained BERT Model Using Labelled Dataset

DGA Detection

In the detection task the objective is to design a classifier to detect whether a domain name is legitimate or malicious successfully. Starting with simple DNN architecture other architectures were designed and tested. Results from these experiments are summarized in Table 1 BERT was proven to be the best performing model, with a precision, recall, and accuracy score of 0.99. The proposed BERT model improves on recent findings published by Yang et al. and Fangali et al. by more than 9 points in absolute terms [5,6].

Furthermore, as seen in Figure 2, the model is highly confident in its classification of test domain names.

Table 1: Detection Results on Test Data. BERT Outperforms all Other Models. P: Precision, R: Recall

| Domain Type | DNN | | CNN | | RNN | | Bi-LSTM | | BERT | | Support |
|--------------|------|------|------|------|------|------|---------|------|-------------|-------------|---------|
| | P | R | P | R | P | R | P | R | P | R | |
| Legit | 0.97 | 0.97 | 0.99 | 0.98 | 0.98 | 0.95 | 0.98 | 0.99 | 0.99 | 0.99 | 56873 |
| AGD | 0.95 | 0.96 | 0.96 | 0.98 | 0.91 | 0.96 | 0.98 | 0.97 | 0.99 | 0.99 | 32992 |
| Macro Avg | 0.96 | 0.97 | 0.97 | 0.98 | 0.94 | 0.95 | 0.98 | 0.98 | 0.99 | 0.99 | 89865 |
| Weighted Avg | 0.97 | 0.97 | 0.98 | 0.98 | 0.95 | 0.95 | 0.98 | 0.98 | 0.99 | 0.99 | 89865 |

DGA Classification

For this task, the problem was defined as classification task in which an output label may belong to either legit or one of 19 DGA classes. Results from test data are shown in Table 2. The classification efficacy of the BERT model is highest for random and dictionary-based algorithms like gozi, matsnu, and supbbox. BERT based classifier outperforms other methods from recent works in Yang L, et al. and Ren F, et al. for both dictionary based and random DGAs [5,6].

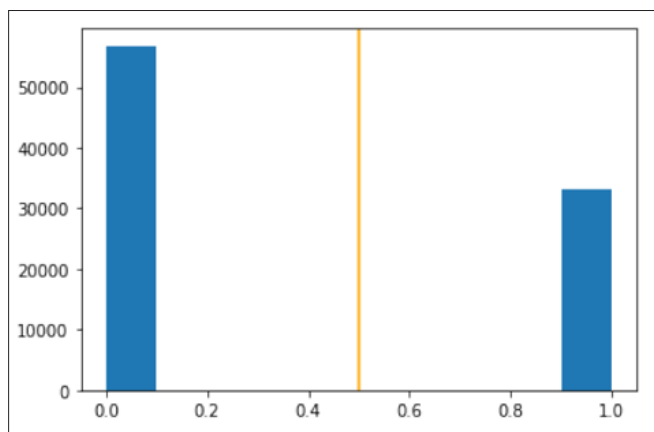


Figure 2: BERT model AGD (1) vs Legitimate(0) Domain Name Detection Confidence. Value of 0.5 was used as Threshold, Shown as Orange Line.

Table 2: Results on Test Data. Dictionary Based DGAs are in Bold. P: Precision, R: Recall

| Domain Type | DNN | | CNN | | RNN | | Bi-LSTM | | BERT | | Support |
|---------------------|------|------|------|------|------|------|---------|------|------|------|---------|
| | P | R | P | R | P | R | P | R | P | R | |
| allureon | 0.86 | 0.94 | 0.68 | 0.88 | 0.88 | 0.97 | 0.9 | 0.98 | 0.87 | 0.95 | 2987 |
| banjori | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2958 |
| cryptolocker | 0.68 | 0.73 | 0.66 | 0.59 | 0.75 | 0.68 | 0.72 | 0.77 | 0.7 | 0.73 | 3007 |
| dyre | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 3032 |
| gozi | 0.8 | 0.98 | 0.94 | 0.97 | 0.84 | 0.98 | 0.96 | 0.97 | 0.99 | 0.97 | 3021 |
| kraken | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2912 |
| legit | 0.94 | 0.94 | 0.97 | 0.97 | 0.96 | 0.96 | 0.98 | 0.97 | 0.99 | 0.99 | 32992 |
| locky | 0.89 | 0.61 | 0.84 | 0.61 | 0.86 | 0.69 | 0.83 | 0.74 | 0.83 | 0.71 | 3024 |
| matsnu | 0.86 | 0.88 | 0.9 | 0.89 | 0.91 | 0.86 | 0.91 | 0.94 | 0.97 | 0.98 | 3028 |
| murofet | 0.98 | 0.99 | 0.99 | 1 | 0.99 | 1 | 0.99 | 1 | 0.99 | 1 | 2995 |
| necurs | 0.97 | 0.79 | 0.97 | 0.81 | 0.97 | 0.82 | 0.98 | 0.83 | 0.99 | 0.83 | 3004 |
| padcrypt | 1 | 1 | 1 | 1 | 0.99 | 1 | 1 | 1 | 1 | 1 | 2975 |
| pushdo | 0.92 | 0.95 | 0.99 | 0.99 | 0.98 | 0.97 | 0.99 | 1 | 0.99 | 0.99 | 3043 |
| pykspa | 0.74 | 0.78 | 0.68 | 0.86 | 0.72 | 0.88 | 0.8 | 0.87 | 0.75 | 0.87 | 2965 |
| qakbot | 0.81 | 0.63 | 0.75 | 0.63 | 0.87 | 0.64 | 0.88 | 0.66 | 0.86 | 0.62 | 3003 |
| ramdo | 0.99 | 1 | 0.99 | 0.99 | 1 | 1 | 1 | 1 | 1 | 1 | 3065 |
| ramnit | 0.77 | 0.73 | 0.75 | 0.62 | 0.8 | 0.78 | 0.79 | 0.84 | 0.75 | 0.82 | 2993 |
| rovnix | 0.89 | 0.79 | 0.96 | 0.93 | 0.9 | 0.83 | 0.95 | 0.96 | 0.97 | 0.99 | 2985 |
| suppobox | 0.89 | 0.99 | 0.94 | 1 | 0.91 | 0.99 | 0.97 | 1 | 0.99 | 1 | 2912 |
| tinba | 0.75 | 0.92 | 0.72 | 0.83 | 0.72 | 0.97 | 0.81 | 0.99 | 0.8 | 0.94 | 2964 |
| <i>macro avg</i> | 0.89 | 0.88 | 0.89 | 0.88 | 0.9 | 0.9 | 0.92 | 0.93 | 0.92 | 0.92 | 89865 |
| <i>weighted avg</i> | 0.9 | 0.9 | 0.91 | 0.91 | 0.92 | 0.92 | 0.94 | 0.94 | 0.94 | 0.94 | 89865 |

Conclusion

Use of algorithmically generated and dictionary based domain names pose a challenge for security experts and network administrators. Current state-of-the-art methods for this task are inadequate for handling this challenge. The proposed model makes use of stacked bi-directional encoder from Transformer architecture. This method is effective and doesn't require any hand-crafted features. Experiments on real word domain names dataset indicate that this model delivers F1 score of 0.99 and more than 0.98 for detection and classification tasks respectively. Proposed mechanism is effective against both random and dictionary based domain generation algorithms.

References

1. Kühner M, Rossow C, Holz T (2014) Paint it black: Evaluating the effectiveness of malware blacklists. International Workshop on Recent Advances in Intrusion Detection 1-21.
2. Antonakakis M, Perdisci R, Nadji Y, Vasiloglou N, Abu-Nimeh S, et al. (2012) From throw-away traffic to bots: detecting the rise of DGA-based malware. 21st {USENIX} Security Symposium ({USENIX} Security 12) 491-506.
3. Woodbridge J, Anderson HS, Ahuja A, Grant D (2016) Predicting domain generation algorithms with long short-term memory networks. arXiv preprint <https://arxiv.org/abs/1611.00791>.
4. Yu B, Gray DL, Pan J, De Cock M, Nascimento AC (2017) Inline dga detection with deep networks. 2017 IEEE International Conference on Data Mining Workshops (ICDMW) 683-692.
5. Yang L, Liu G, Dai Y, Wang J, Zhai J (2020) Detecting stealthy domain generation algorithms using heterogeneous deep neural network framework. Ieee Access 8: 82876-82889.
6. Ren F, Jiang Z, Wang X, Liu J (2020) A DGA domain names detection modeling method based on integrating an attention mechanism and deep neural network. Cybersecurity <https://cybersecurity.springeropen.com/articles/10.1186/s42400-020-00046-6>.
7. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, et al. (2017) Attention is all you need. In: Advances in neural information processing systems 5998-6008.
8. Devlin J, Chang MW, Lee K, Toutanova K (2018) Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint <https://arxiv.org/abs/1810.04805>.
9. Zago M, Perez MG, Perez GM (2020) Umudga: A dataset for profiling algorithmically generated domain names in botnet detection. Data in brief <https://www.sciencedirect.com/science/article/pii/S2352340920302948>.

Copyright: ©2024 Santosh Kumar Kande. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.