

# Towards Domain Generating Algorithm based Malicious Domains Detection

Md. Ahsan Ayub and Steven Smith  
Advisor: Dr. Ambareen Siraj  
Department of Computer Science

## ABSTRACT

- A machine learning approach for effective detection of malicious Domain Generating Algorithm (DGA) based Domains used by botnets and other malware for evasion.
- Makes use of two feature extraction methods, Bag of Words and Word2Vec for text processing.
- Considers binary detection and multiclass classification for 84 different DGA families, the largest study of DGA domain detection to date.

## BACKGROUND

- DGAs are used to dynamically produce a large set of domains to evade blacklisting and reverse engineering.
- Two types of DGAs: Traditional DGAs & Dictionary-based DGAs.
- DGAs are primarily used by botnets to aid in performing cyberattacks such as DDOS, and in sending spam and phishing emails.

## DATASET

- DGArchive: DGA domains labelled by DGA family (84 total families).
- Majestic Million: Top 1 Million Most Visited domains used as benign Domain Names
- Dataset split into 70% training, 20% testing, and 10% validation.

## DETECTION METHODOLOGY

- Two techniques compared using the Bag of Words (BoW) Bigram model, and the Word2Vec model.
- Considered NXDomain and VirusTotal Scan Results for Classification.
- Detection with the BoW Model:
  - Bigram (2-Gram) Model used to capture context of two word combinations in domains.
  - Logistic Regression, Decision Tree, and Artificial Neural Network (ANN) considered.
- Detection with the Word2Vec Model:
  - Long Short Term Memory (LSTM) Network used to capture temporal relationships among tokens in a sequence.

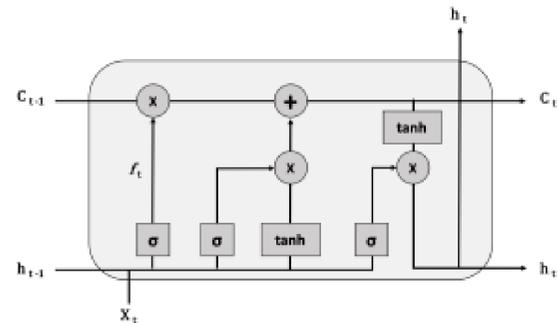


Fig. 1: A Unit Cell of the Long Short Term Memory (LSTM) Network

## RESULTS

Model	Accuracy	Precision	Recall	$F_1$
Logistic Regression	0.9816	0.9911	0.9993	0.9833
Decision Tree	0.9965	0.9965	0.9988	0.9941
<b>ANN</b>	<b>0.9979</b>	<b>0.9979</b>	<b>0.9979</b>	<b>0.998</b>
LSTM	0.995	0.9949	0.9958	0.994

Fig. 2: Performance of Each Model for Binary Classification

- Binary Classification/Detection
  - ANN with Bigram BoW model proves to be the highest performing technique – With over 99% accuracy, precision,  $F_1$ , and Recall Scores.
  - Best classification results to date as seen in Fig. 3.

Research Work	No. of DGAs	Method	Classification	Precision	Recall	$F_1$
Our Study	84	ANN	Binary	<b>0.9979</b>	<b>0.9979</b>	<b>0.998</b>
			Multiclass	0.9358	0.9358	0.9358
Woodbridge et al. [92]	30	LSTM	Binary	0.9942	0.9937	0.9906
			Multiclass	<b>0.963</b>	<b>0.97</b>	<b>0.963</b>
Lison et al. [55]	56	RNN	Binary	0.972	0.97	0.971
			Multiclass	0.891	0.892	0.887
Tran et al. [88]	37	LSTM	Binary	0.9842	0.9842	0.9842
			Multiclass	0.8728	0.8775	0.8751

Fig. 3: Performance Compared to Previous Work

- Multiclass Classification
  - Effective classification for 69 out of 84 DGA families.
  - Considers the most families out of any study.
  - Average performance lowered due to 10 families with less than 300 samples, and 5 with similar randomness to benign samples.

## CONCLUSION

- Achieved the best results to date for Binary Classification with the combination of Bigram BoW and ANN.
- Effective multiclass classification for a larger set of families than previous studies through the use of Bigram BoW and ANN techniques.

## ACKNOWLEDGEMENT

The work reported in this poster has been fully supported by Cybersecurity Education, Research & Outreach Center (CEROC) at Tennessee Tech.