

Goal

To identify potential malware network traffic and classify Ransomware family via machine learning.

- Malware traffic identification
- Ransomware family classification

Feature Extraction

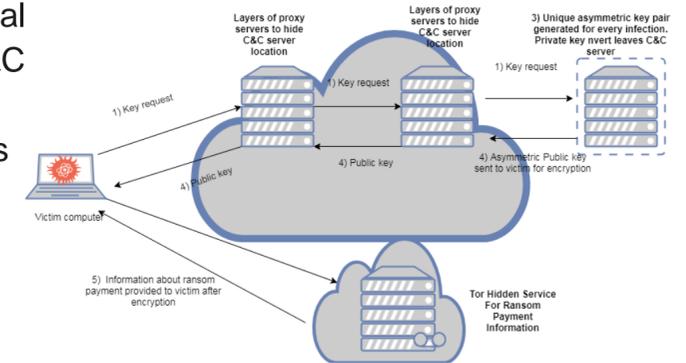
- Connection Records : port, IP, and protocol information in traffic flows from the transport layer.
- Network Packet Payload : extract all the payloads from each packet where the length of each packet ranges from 0 to 1,500 bytes.
- Flow Behavior : inter-arrival time of each packet and Markov transition matrix to represent the flow behavior and document the before-and-after relationship.

Dataset Overview

Class	Number of Flows	Size
Benign	246,015	560.2 MB
Bot	99	6.5 MB
Exploit	349	32.5 MB
Trojan	3,085	18.1 MB
Malspam	3,612	142.1 MB
Cryptomix	90	2.0 MB
Locky	229	9.3 MB
CrypMic	390	14.3 MB
Teslacrypt	755	26.5 MB
CryptXXX	1,259	44.7 MB
Cryptowall	2,864	34.7 MB
Cerber	23,260	23.5 MB
Total	282,007	914.4 MB

Ransomware Overview

- Victim makes initial key request to C&C server
- C&C server returns encryption Public key
- Tor Hidden Service For Ransom Payment Information



Quantity Dependent Backpropagation (QDBP)

To mitigate imbalance data issue, we introduce a vector F into backpropagation and propose a QDBP algorithm which takes the disparity between classes into consideration and shows different sensitivities toward different classes. The mathematical formula is given by:

$$\theta_i^{l+} = \theta_i^l - \eta \cdot F \cdot \nabla Loss$$

$$F = \left[\frac{c_1}{n_1}, \frac{c_2}{n_2}, \dots, \frac{c_N}{n_N} \right]$$

$$\nabla Loss = \left[\frac{\partial Loss_1}{\partial \theta_i}, \frac{\partial Loss_2}{\partial \theta_i}, \dots, \frac{\partial Loss_N}{\partial \theta_i} \right]^T$$

Malware traffic identification Results

Method	Accuracy	Precision
DNN + Backpropagation	59.08%	8.33%
DNN + Oversampling (10000 samples/class)	85.18%	65.9%
DNN + Undersampling (45 samples/class)	68.89%	49.45%
DNN + Incremental Learning	78.84%	71.23%
DNN + QDBP	84.56%	62.3%
SVM (RBF)	83.87%	38.8%
Random Forest	98.9%	68.25%
TSDNN + QDBP	99.63%	85.4%

Key Flow Features

- Packet inter-arrival times
- Network Packet Payload
- Markov transition matrix

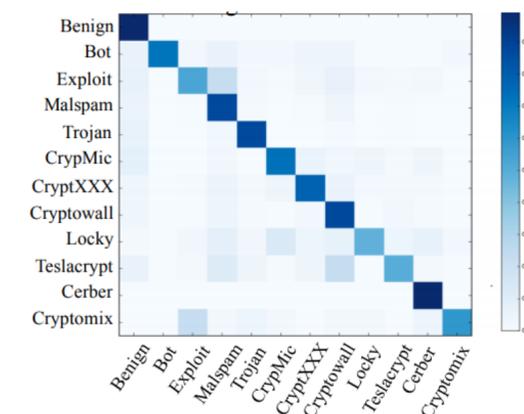
Ransomware family classification

Goal

- Minimize false positives

Performance

- Precision: 0.964
- Recall: 0.943
- F1 Score: 0.952
- AUC of ROC: 0.971

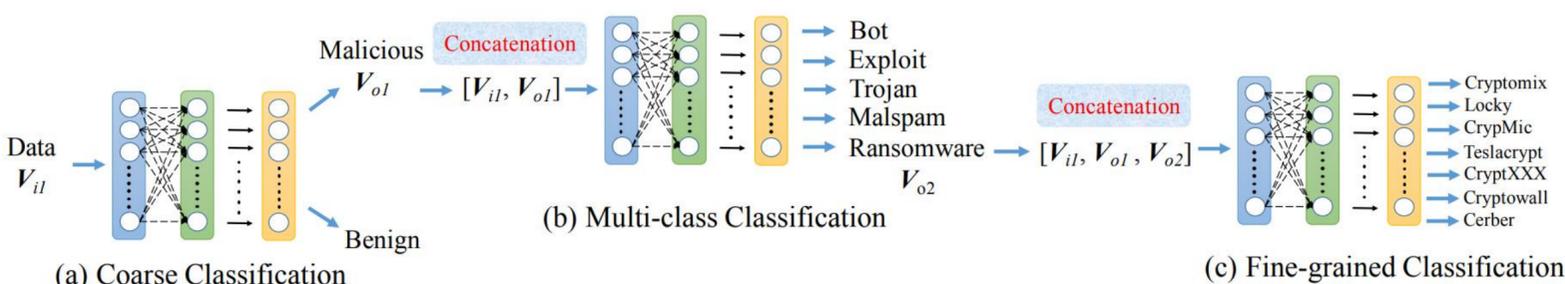


Key Flow Features

- Packet inter-arrival times
- Network Packet Payload
- Http Request

Tree-Shaped Deep Neural Network (TSDNN)

To mitigate the imbalanced data issue, we propose an end-to-end trainable TSDNN model which classifies the data layer by layer. We define each model in TSDNN as a nodal network and the links between nodal networks are called bridges. We adopt cross entropy as the loss function and apply QDBP to each nodal network to optimize the performance.



Discussion

The experimental results show that our model is able to accurately detect the potential malware which also justifies that behavior-oriented approach is better than signature-oriented methods when detecting malware.