

## PROBLEM

Cyber attacks are continuously evolving and becoming more sophisticated and difficult to detect.

## QUESTIONS

How do deep learning models detect network intrusions?

Can they be bypassed by adjusted attacks?

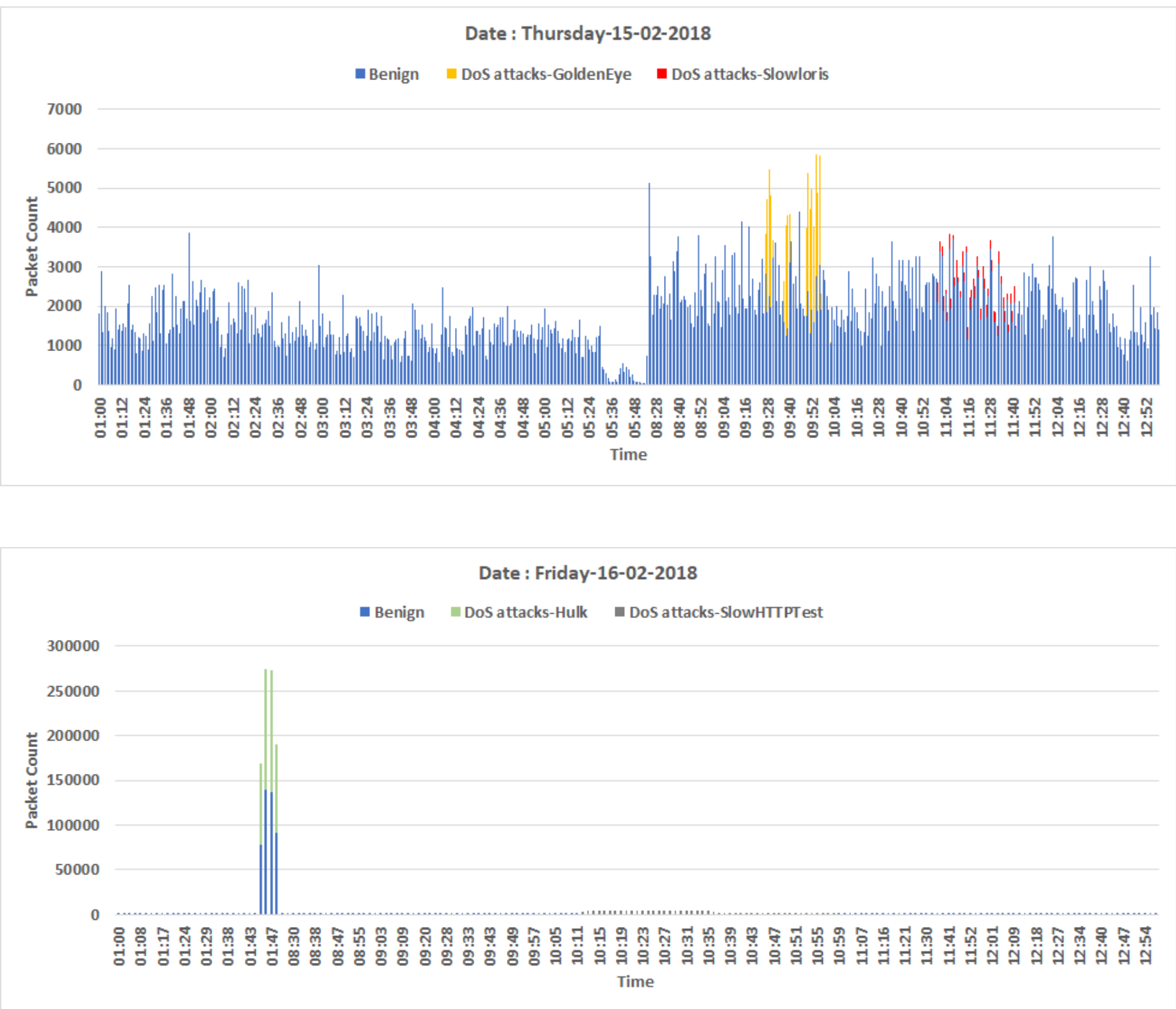
If so, what approaches are less easily fooled?

## GOAL

To experiment with various adversarial attacks, custom intrusion sets, and various neural network architectures, employing Communications Security Establishment (CSE) & the Canadian Institute of Cybersecurity (CIC) intrusion data set of 2018 that reflects current network trends.

## DATA SET

75 Traffic Features of 6.1 Million Flows in total

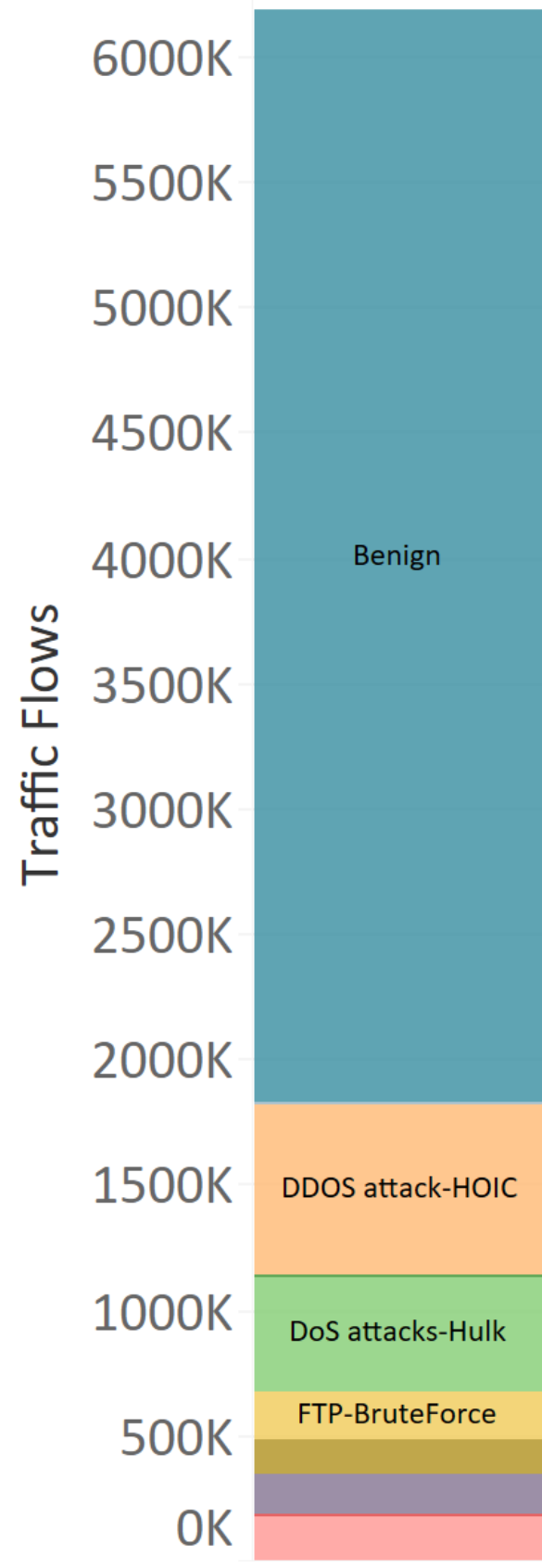


CSE-CIC-IDS2018

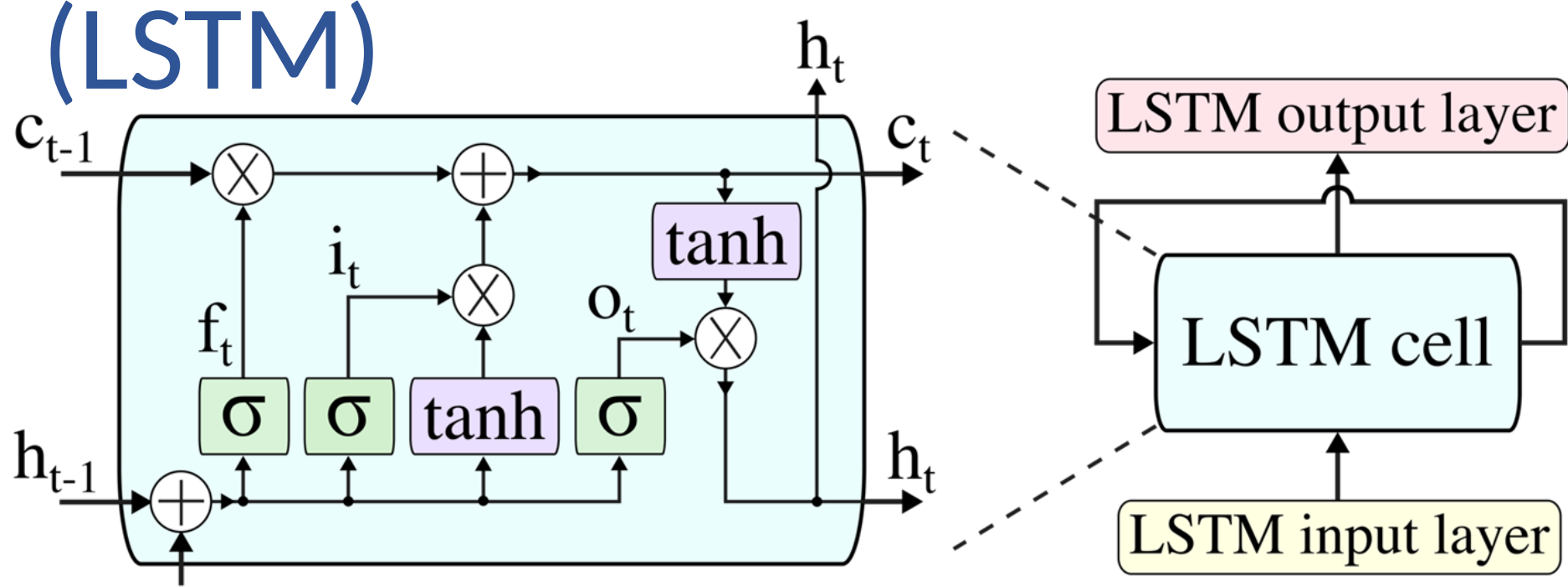
Testbed: Attacker-network with 50 terminals, and victim-network, implemented as a Local Area Network (LAN) with 420 terminals and 30 servers divided into 5 subnets [1,2].

Benign data: background traffic based on user non-malicious behavior when employing: HTTP, HTTPS, SMTP, POP3, IMAP, SSH, and FTP.

- Label
- Benign
  - Brute Force -Web
  - Brute Force -XSS
  - DDoS attack-HOIC
  - DDoS attack-LOIC-UDP
  - DoS attacks-Hulk
  - FTP-BruteForce
  - Infiltration
  - SQL Injection
  - SSH-Bruteforce



## Long Short Term Memory (LSTM)



The outputs of the forget gate  $f_t$ , the input gate  $i_t$ , and the output gate  $o_t$  at time  $t$  are:  
 $f_t = \sigma(W_{ff}x_t + b_{ff} + U_{ff}h_{t-1} + b_{hf})$   
 $i_t = \sigma(W_{if}x_t + b_{if} + U_{if}h_{t-1} + b_{hi})$   
 $o_t = \sigma(W_{of}x_t + b_{of} + U_{of}h_{t-1} + b_{ho})$

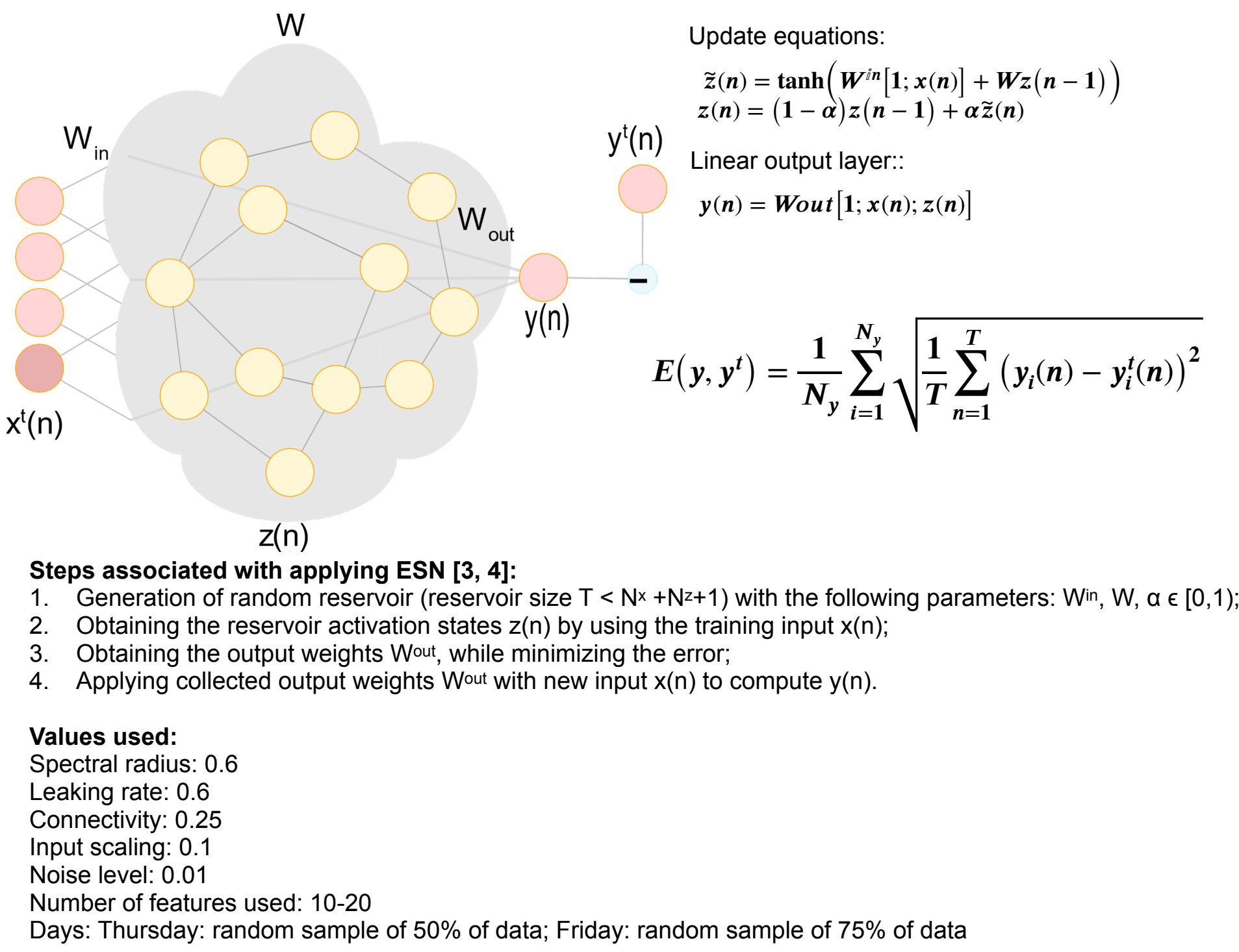
The cell state  $c_t$  is calculated as:

$$c_t = f_t c_{t-1} + i_t \tanh(W_{ic}x_t + b_{ic} + U_{ic}h_{t-1} + b_{hc})$$

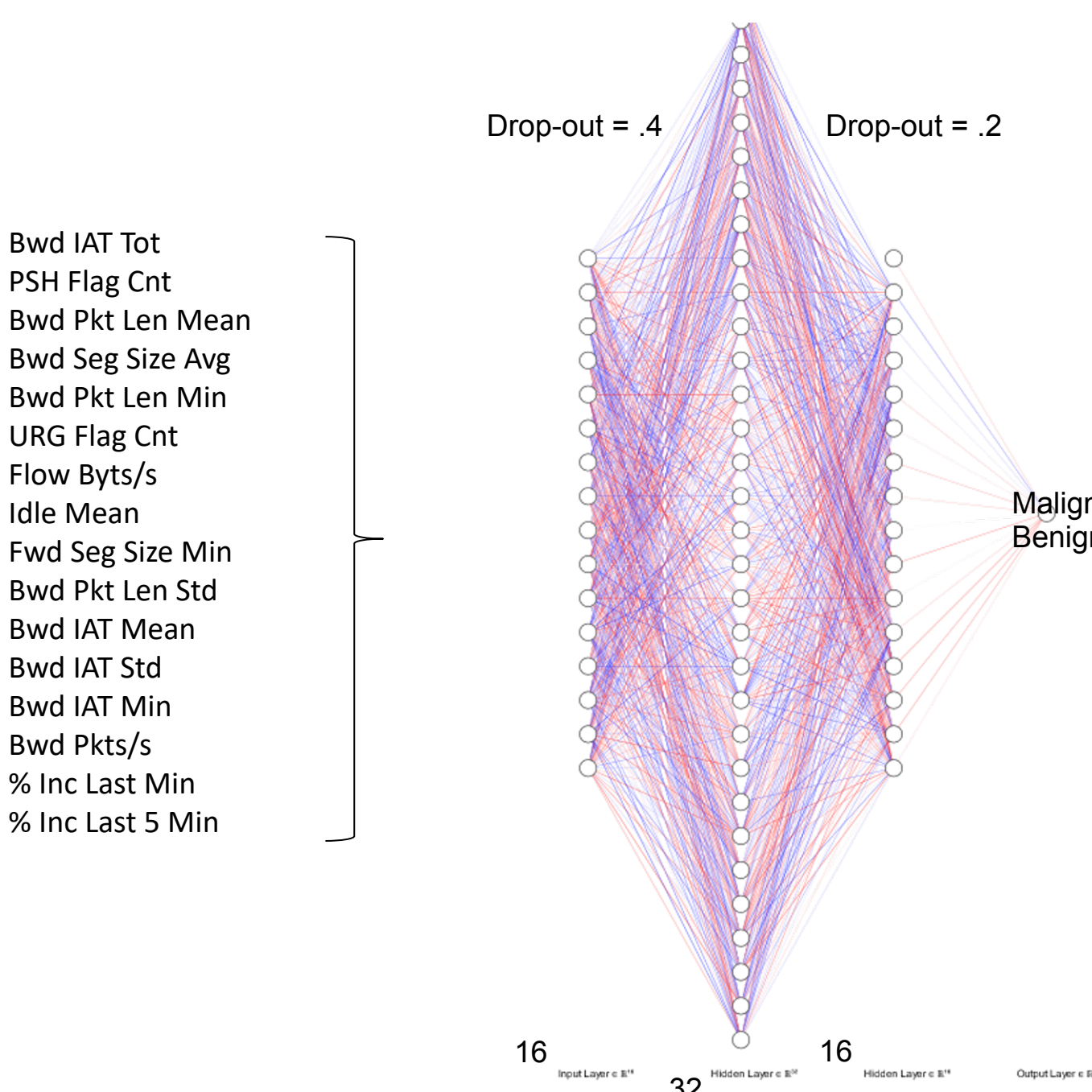
The output of the LSTM cell is:

$$h_t = o_t \tanh(c_t)$$

## Echo State Networks (ESN)

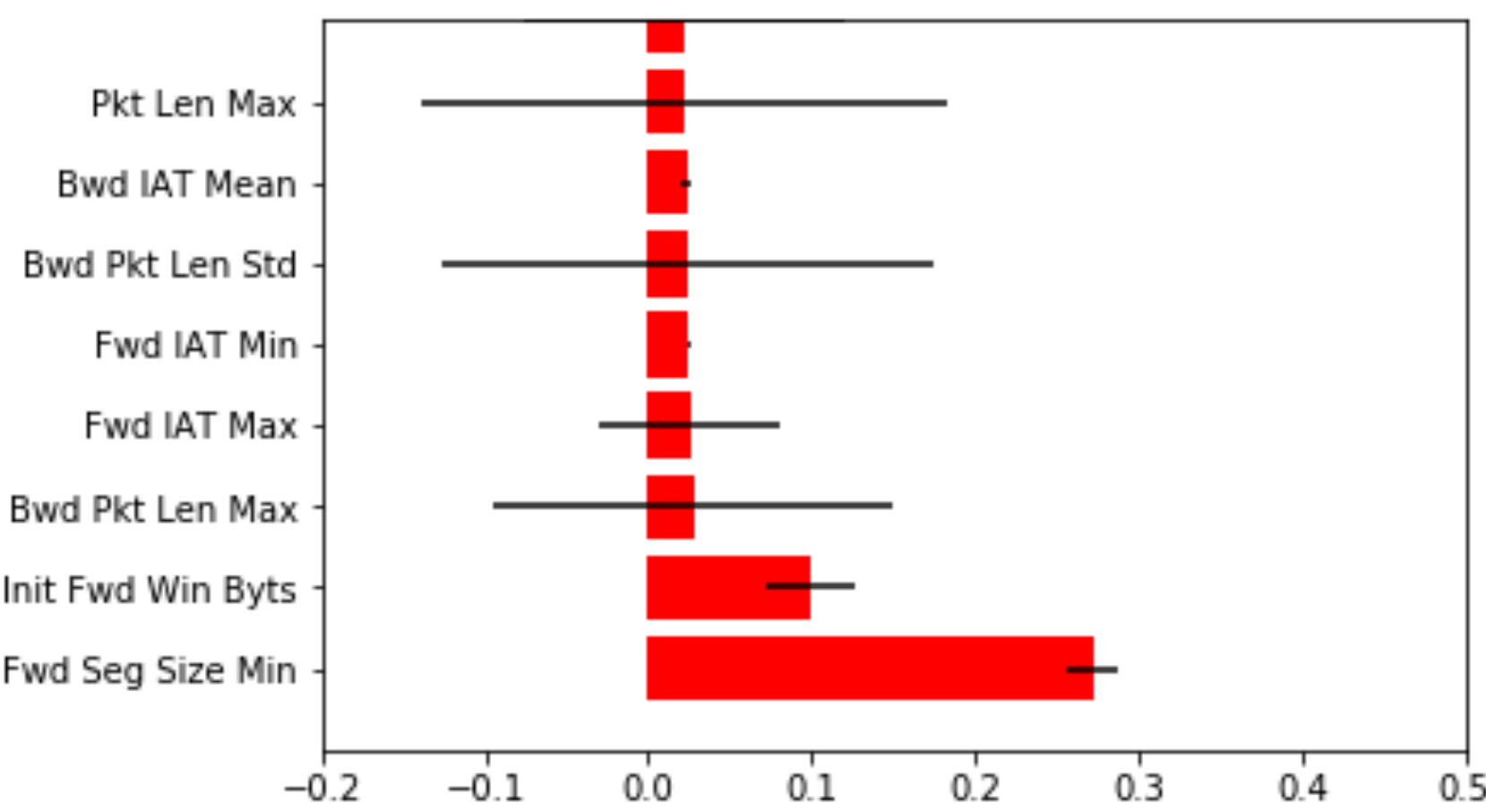


## Feed-Forward Neural Networks (FFNN)

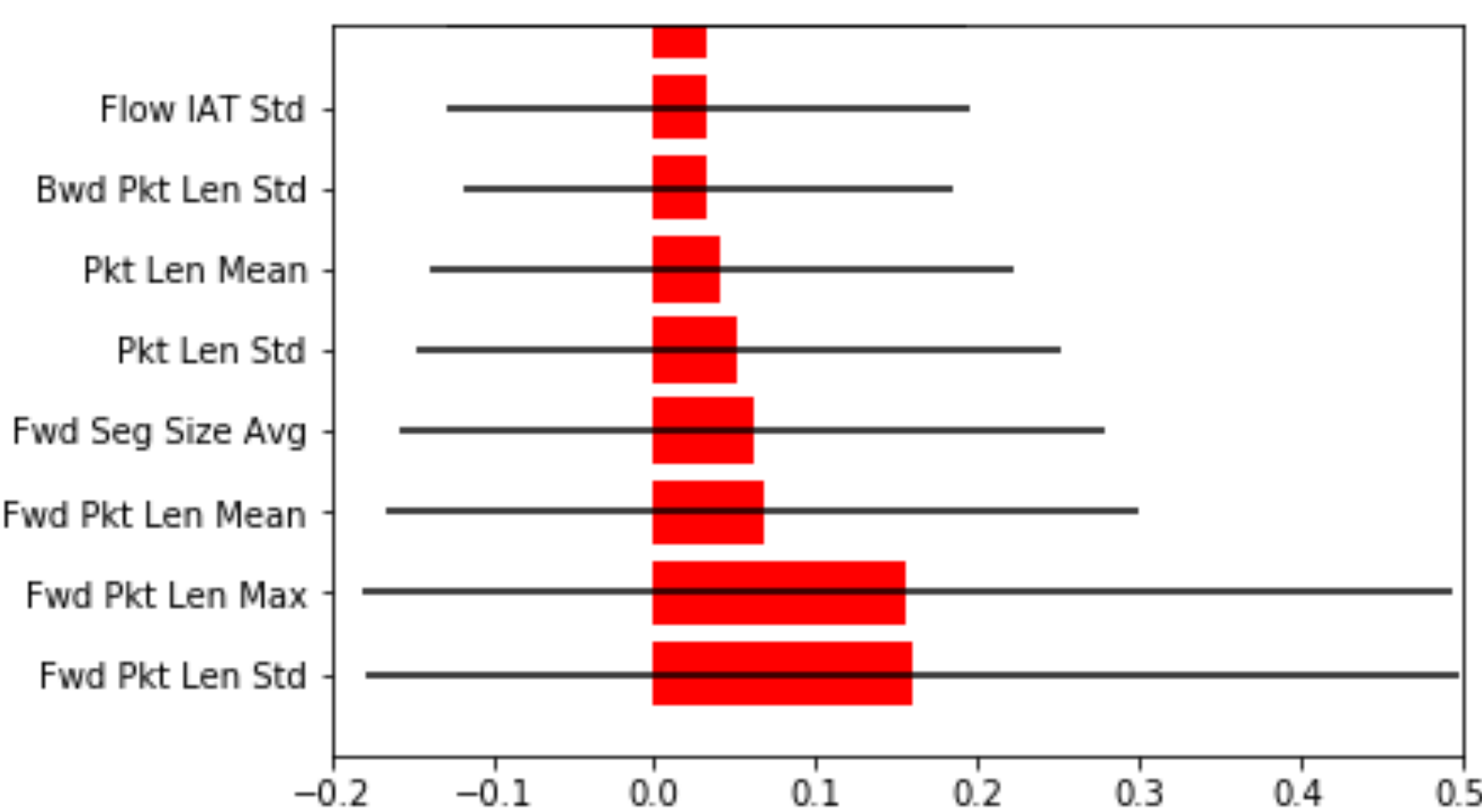


## EXPERIMENTS (Feature selection and adversarial landscape)

Importance of Traffic Features for detecting GoldenEye/Slowloris Attacks



Importance of Traffic Features for detecting Hulk/Slow HTTP Attacks

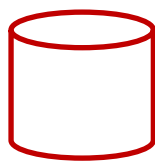


Feature Engineering/Selection

Trained Models on 2018 attacks



Additional Test on Modified 2019 Attack



## RESULTS

LSTM

LSTM Test Set Hulk/Slow HTTP Attack: 99.5%

|                  | Actual |         |
|------------------|--------|---------|
|                  | Benign | Attack  |
| Predicted Benign | 88,450 |         |
| Predicted Attack | 947    | 120,318 |

Test Set GoldenEye/Slowloris Attacks: 95.0%

|                  | Actual  |        |
|------------------|---------|--------|
|                  | Benign  | Attack |
| Predicted Benign | 199,226 | 10,489 |
| Predicted Attack |         |        |

LSTM New 2019 Hulk Attack: 97.1%

|                  | Actual |        |
|------------------|--------|--------|
|                  | Benign | Attack |
| Predicted Benign | 9      |        |
| Predicted Attack | 249    | 8,450  |

ESN

ESN Test Set Hulk/Slow HTTP Attack: 98.9%

|                  | Actual |         |
|------------------|--------|---------|
|                  | Benign | Attack  |
| Predicted Benign | 82,573 | 1,411   |
| Predicted Attack | 847    | 111,777 |

Test Set GoldenEye/Slowloris Attacks: 93.9%

|                  | Actual  |        |
|------------------|---------|--------|
|                  | Benign  | Attack |
| Predicted Benign | 122,927 | 6,497  |
| Predicted Attack | 1,561   | 87     |

ESN New 2019 Hulk Attack: 97.2%

|                  | Actual |        |
|------------------|--------|--------|
|                  | Benign | Attack |
| Predicted Benign | 11     | 0      |
| Predicted Attack | 247    | 8,450  |

FFNN

FFNN Test Set Hulk/Slow HTTP Attack: 99.8%

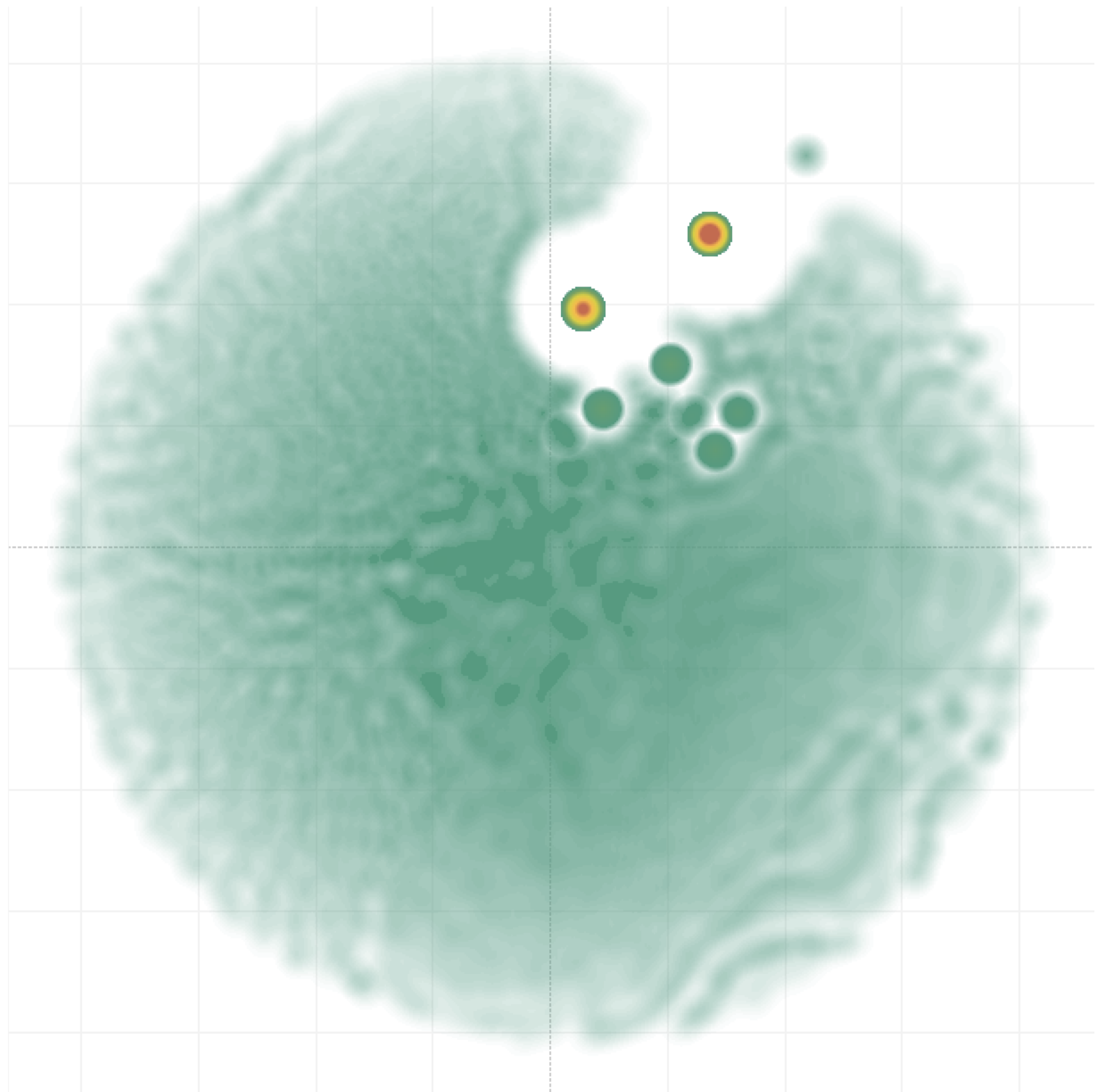
|                  | Actual |         |
|------------------|--------|---------|
|                  | Benign | Attack  |
| Predicted Benign | 89,265 | 337     |
| Predicted Attack | -      | 120,113 |

Test Set GoldenEye/Slowloris Attacks: 99.7%

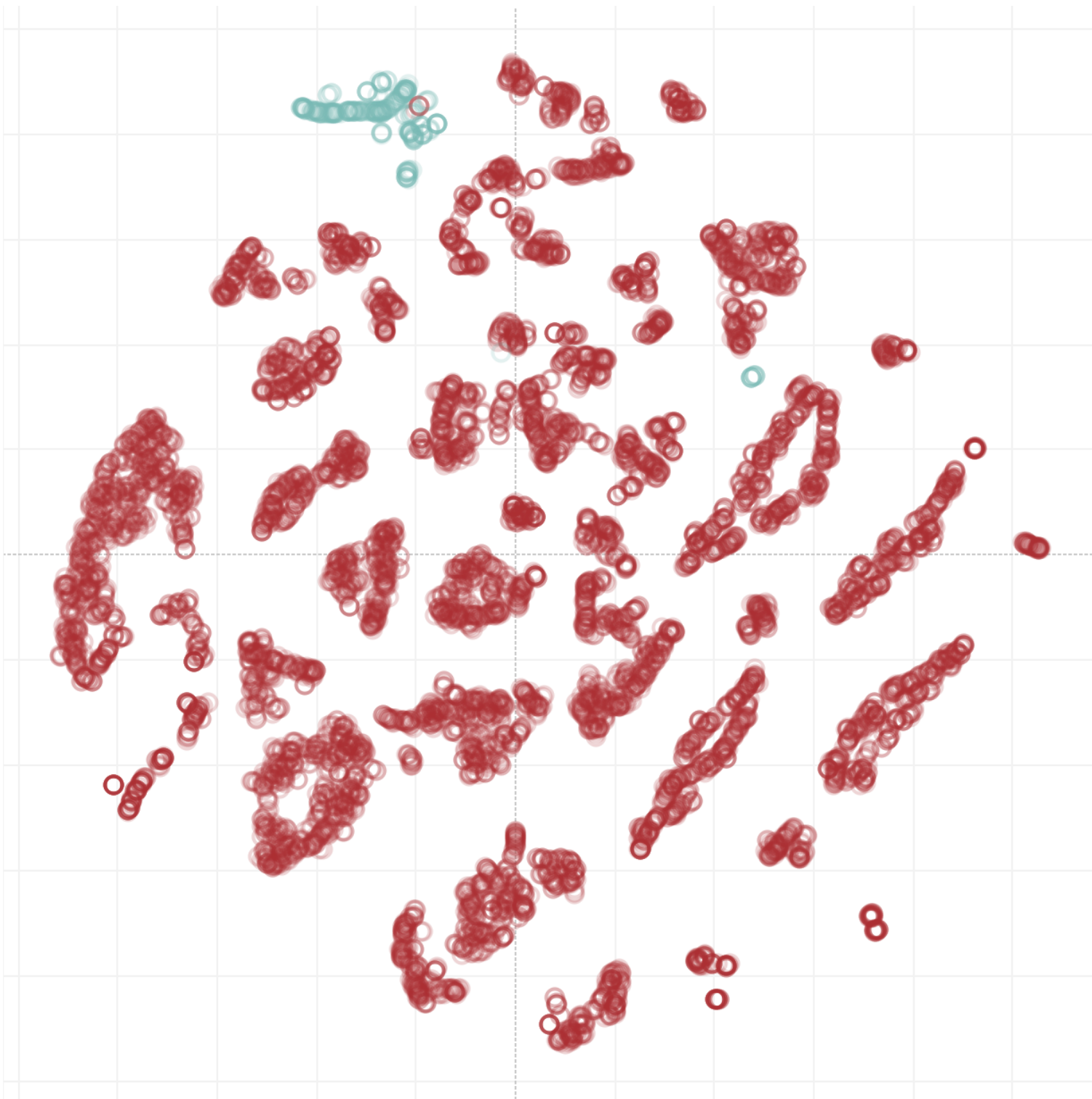
|                  | Actual  |        |
|------------------|---------|--------|
|                  | Benign  | Attack |
| Predicted Benign | 199,005 | 145    |
| Predicted Attack | 572     | 9,993  |

FFNN New 2019 Hulk Attack: 3.2%

|                  | Actual |        |
|------------------|--------|--------|
|                  | Benign | Attack |
| Predicted Benign | 258    | 8,432  |
| Predicted Attack | -      | 18     |



T-SNE of a Hulk DOS Attack versus the 99% Accurate FFNN Model



T-SNE of a Modified Hulk DOS Attack

The red marks are all the examples that the same FFNN Model pictured above 'missed'

## References:

- [1] Intrusion Detection Evaluation Dataset (CICIDS2017) [Online]. Available: <https://www.unb.ca/cic/datasets/ids-2017.html>. Accessed: Oct. 28, 2019.
- [2] CICFlowMeter [Online]. Available: <http://netflowmeter.ca/netflowmeter.html>. Accessed: Oct. 28, 2019.
- [3] M. Lukosevicius, "A Practical Guide to Applying Echo State Networks," in *Neural Networks: Tricks of the Trade*. Springer, 2012, pp. 659-686.
- [4] H. Jaeger, "The 'echo state' approach to analysing and training recurrent neural networks," Tech. Rep. GMD Rep. 148, German Nat. Res. Center for Inf. Technol., 2001.
- [5] Z. Li, P. Batta, and Lj. Trajkovic, "Comparison of machine learning algorithms for detection of network intrusions," in Proc. IEEE International Conference on Systems, Man, and Cybernetics (SMC 2018), Miyazaki, Japan, Oct. 2018, pp. 4248-4253.