



Abdullah Al-Dujaili, Alex Huang, Erik Hemberg, Una-May O'Reilly

### Malware detectors have been improved $\bullet$ and automated with the help of machine learning (ML).

- Unfortunately adversaries aiming to thwart lacksquarethem are also aided by ML (see Fig. 1).
- Our project explores the malware arms  $\bullet$ race of artificially intelligent competitors for ways detectors can gain the upper hand.

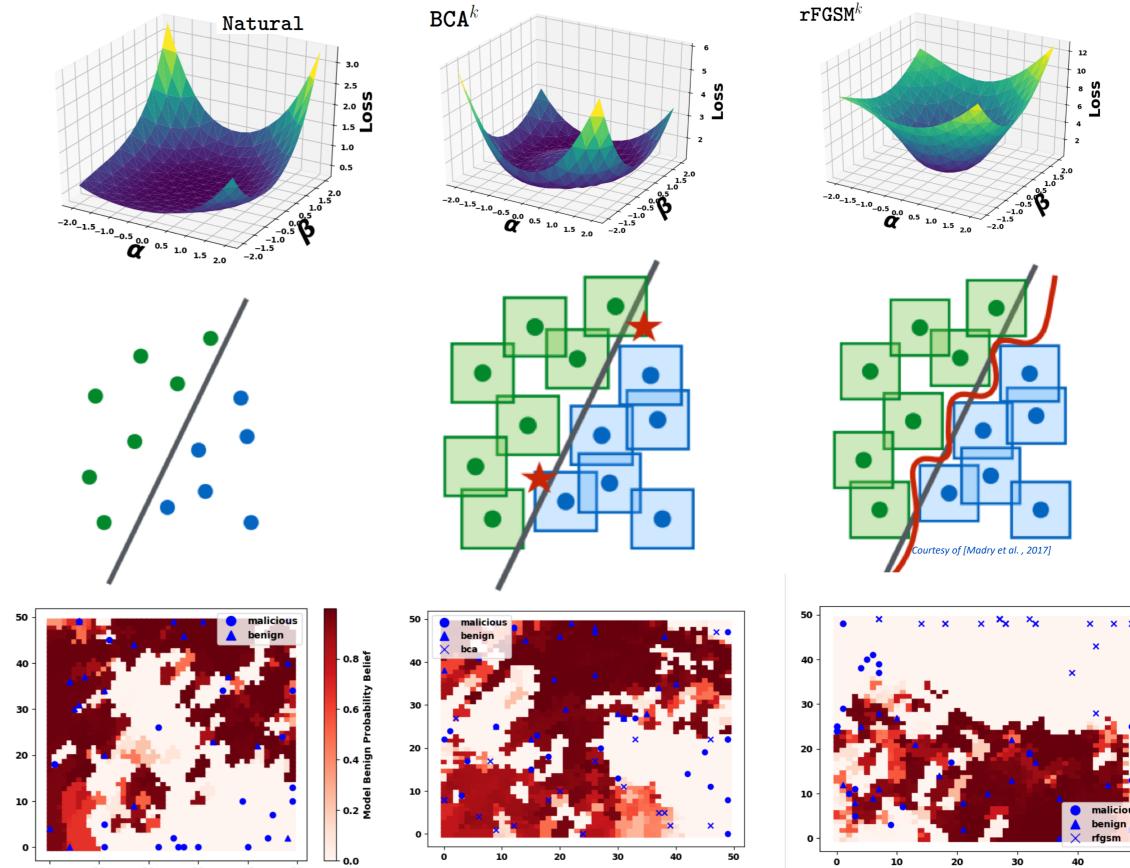
# **Project Overview** func in binary.imported\_functions • • 1 file\_1.exe (PE)

Fig. 1. An appropriate feature vector of a PE can be built based on the PE's imported functions. A malware author can import additional functions without affecting the malicious functionality giving rise to a set of adversarial versions of the same PE. These adversarial versions can be found with the help of an ML algorithm.

# **Current Progress**

 $\mathbf{x}^{(1)} \mathbf{x}^{(2)} \mathbf{x}^{(3)} \cdots \mathbf{x}^{(r)}$ 

 $\mathbf{x}_{ada}^{(1)}$   $\mathbf{x}_{ada}^{(2)}$   $\mathbf{x}_{ada}^{(3)}$  ...  $\mathbf{x}_{ada}^{(r)}$ 



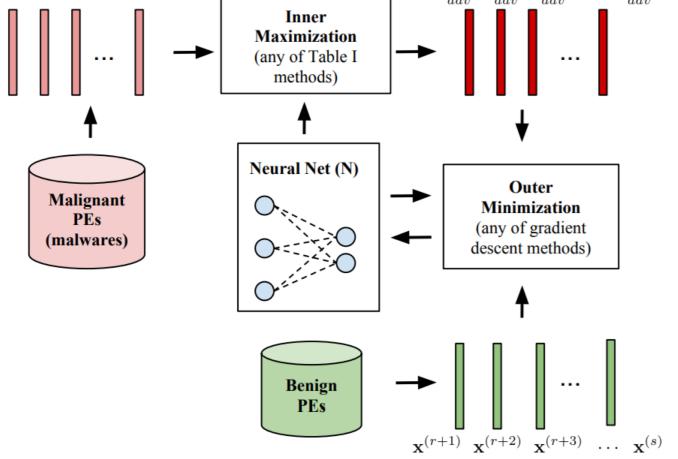


Fig. SLEIPNIR: saddle-point 2. а formulation for hardening ML models with binary feature space.

Table I. Performance Metrics

| Model              | Accuracy | FPR  | FNR  | $\bar{\mathcal{N}}_{\mathrm{BS}}$ |
|--------------------|----------|------|------|-----------------------------------|
| Natural            | 91.9     | 8.2  | 8.1  | 1.0                               |
| $dFGSM^k$          | +0.1     | +1.4 | -1.7 | +1.6                              |
| ${\tt rFGSM}^k$    | -0.6     | +3.6 | -2.4 | +3.0                              |
| $\mathtt{BGA}^k$   | +0.2     | +0.0 | -0.5 | +2.5                              |
| $BCA^k$            | -0.3     | +0.9 | -0.5 | +0.0                              |
| [Grosse et al.'17] | -1.1     | -3.9 | +5.9 | +0.6                              |

| Table | e II. | Evasion | Rates |
|-------|-------|---------|-------|
|-------|-------|---------|-------|

| Model              | Adversary  |           |                 |            |         |                  |  |
|--------------------|------------|-----------|-----------------|------------|---------|------------------|--|
|                    | Natural    | $dFGSM^k$ | ${\tt rFGSM}^k$ | $BGA^k$    | $BCA^k$ | [Grosse et al.'1 |  |
| Natural            | 8.1        | 99.7      | 99.7            | 99.7       | 41.7    | 99.7             |  |
| $dFGSM^k$          | 6.4        | 6.4       | 21.1            | 7.3        | 27.4    | 99.2             |  |
| ${\tt rFGSM}^k$    | <b>5.7</b> | 7.0       | <b>5.9</b>      | <b>5.9</b> | 6.8     | 35.0             |  |
| $BGA^k$            | 7.6        | 39.6      | 17.8            | 7.6        | 10.9    | 68.4             |  |
| $\mathbf{BCA}^k$   | 7.6        | 99.5      | 99.5            | 91.8       | 7.9     | 98.6             |  |
| [Grosse et al.'17] | 14.0       | 69.3      | 69.3            | 37.5       | 14.1    | 15.6             |  |

Fig. 3. In contrast to standard generalization, the decision landscape (input space) has a stronger association with the hardened model's robust generalization, compared to the geometry of the loss landscape (parameter space).

#### **Relevant Papers**

1. Al-Dujaili, et al. "Adversarial Deep Learning for Robust Detection of Binary Encoded

Malware." IEEE SPW, 2018.

2. Huang et al. "On Visual Hallmarks of Robustness to Adversarial Malware." IJCAI-IREDLIA, 2018..

### **Open Questions & Future Vision**

- Static Analysis vs. Dynamic Analysis
- **Representation Learning: Abstract** Syntax Trees (ASTs)
- Interpreted languages: PowerShell Scripts

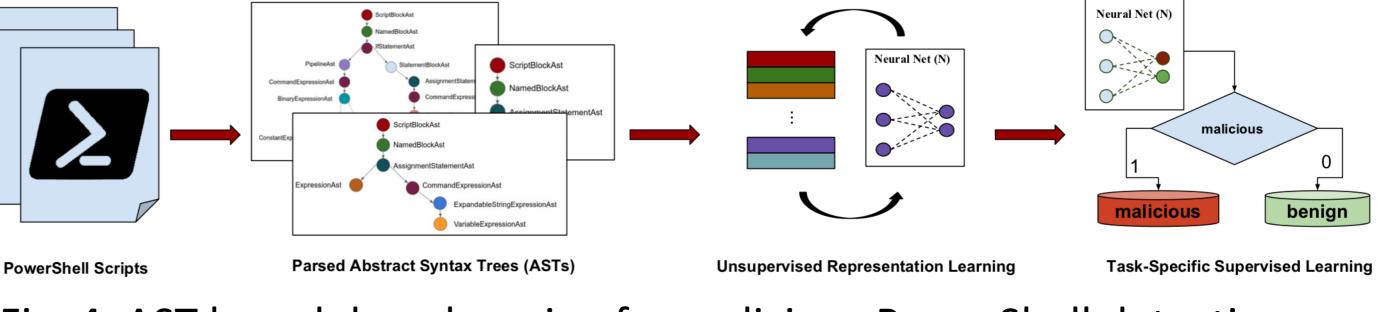


Fig. 4. AST-based deep learning for malicious PowerShell detection.

## **#AI Research Week**

hosted by MIT-IBM Watson AI Lab

**Detailed Contact Information** (corresponding author): Abdullah Al-Dujaili, aldujail@mit.edu