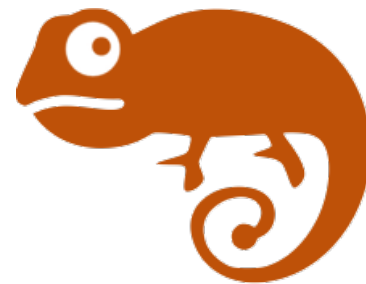# Automatically Evading Classifiers
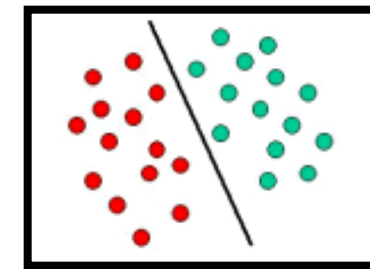## A Case Study on PDF Malware Classifiers

Weilin Xu

David Evans

Yanjun Qi

University of Virginia

# Machine Learning is Solving Our Problems

Spam          IDS          Fake Accounts          Malware          …

**kaggle** / **Microsoft**

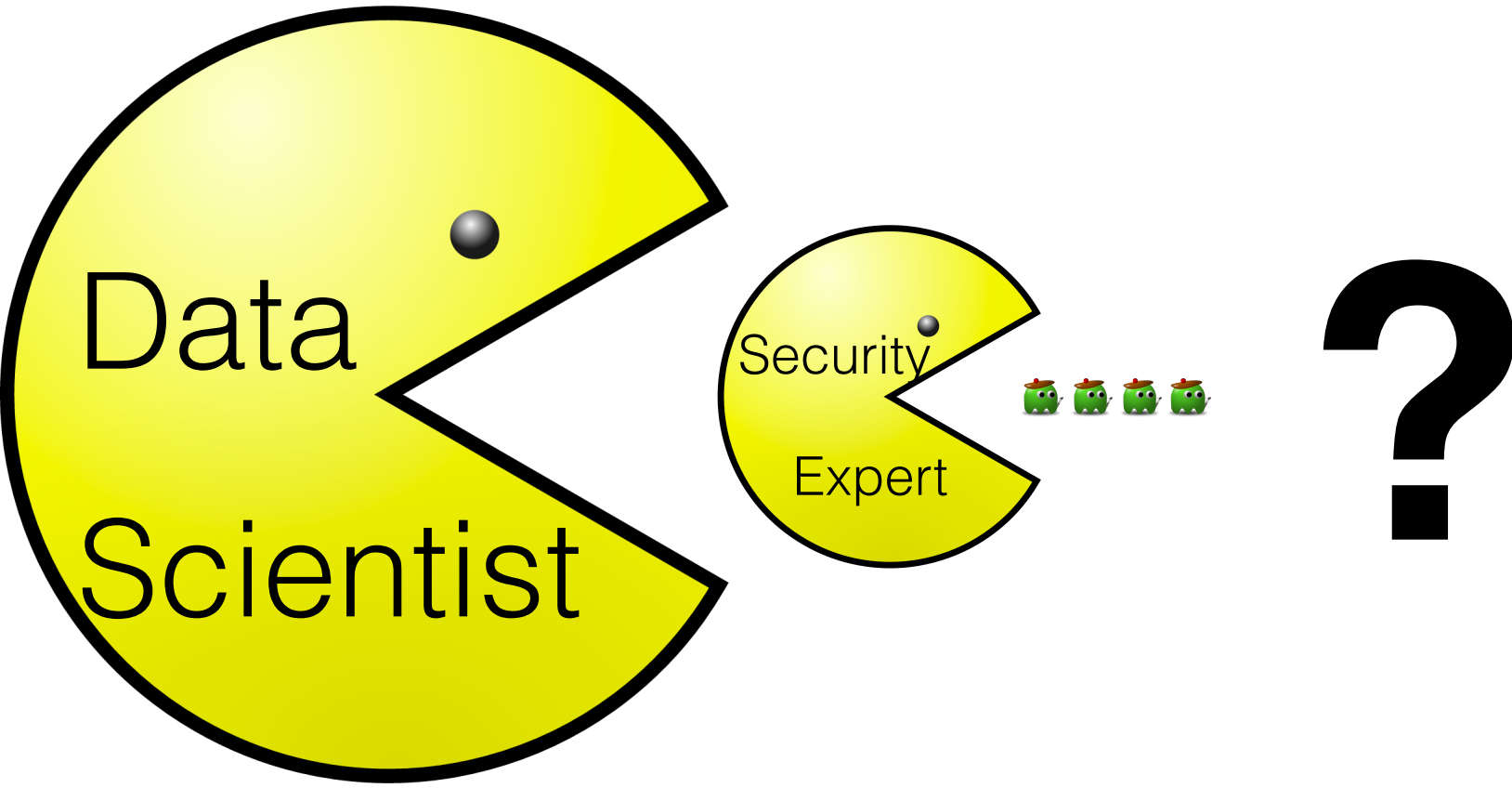# Microsoft Malware Classification Challenge (BIG 2015)

Tue 3 Feb 2015 – Fri 17 Apr 2015 (10 months ago)

| # | Δrank | Team Name  * in the money | Score | Entries | Last Submission UTC (Best – Last Submission) |
|---|-------|---------------------------|-------|---------|----------------------------------------------|
| 1 | ↑5 | ○ say NOOOOO to overfitttttting 👥 *<br>• Little Boat<br>• rcarson<br>• Xueer Chen | 0.002833228 | 268 | Fri, 17 Apr 2015 23:21:56 |
| 2 | ↑7 | Marios & Gert 👥 * | 0.003240502 | 80 | Fri, 17 Apr 2015 12:13:53 (-25.4h) |
| 3 | ↑11 | ○ Mikhail & Dmitry & Stanislav 👥 * | 0.003969846 | 71 | Fri, 17 Apr 2015 23:54:08 |
| 4 | ↑13 | Ivica Jovic | 0.004470816 | 11 | Fri, 17 Apr 2015 23:53:38 (-0.2h) |
| 5 | ↑8 | Octo Guys 👥 | 0.005191324 | 37 | Fri, 17 Apr 2015 23:54:57 (-1.5h) |
| 6 | ↑12 | ○ Oleksandr Lysenko | 0.005335339 | 51 | Fri, 17 Apr 2015 20:26:27 (-12.5h) |

# Machine Learning is Eating the World

# Machine Learning is Eating the World

Data
Scientist

Security
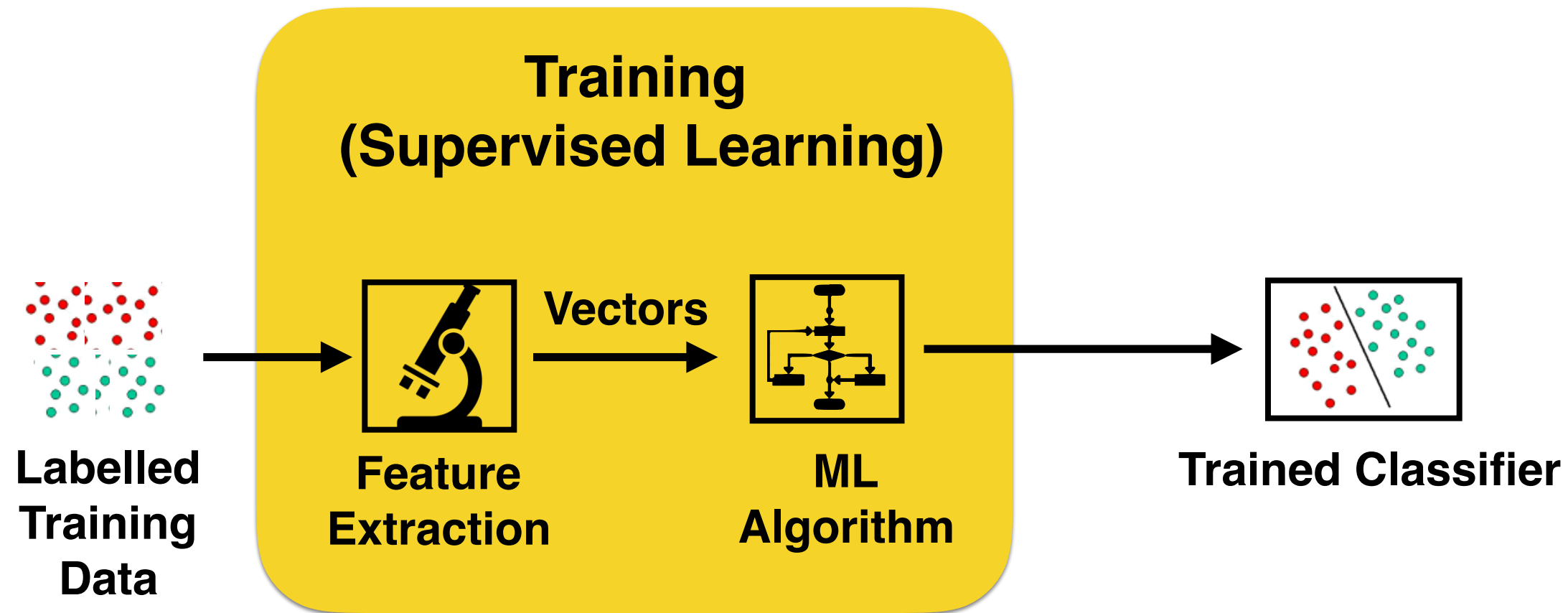Expert

**No!
Security is different.**

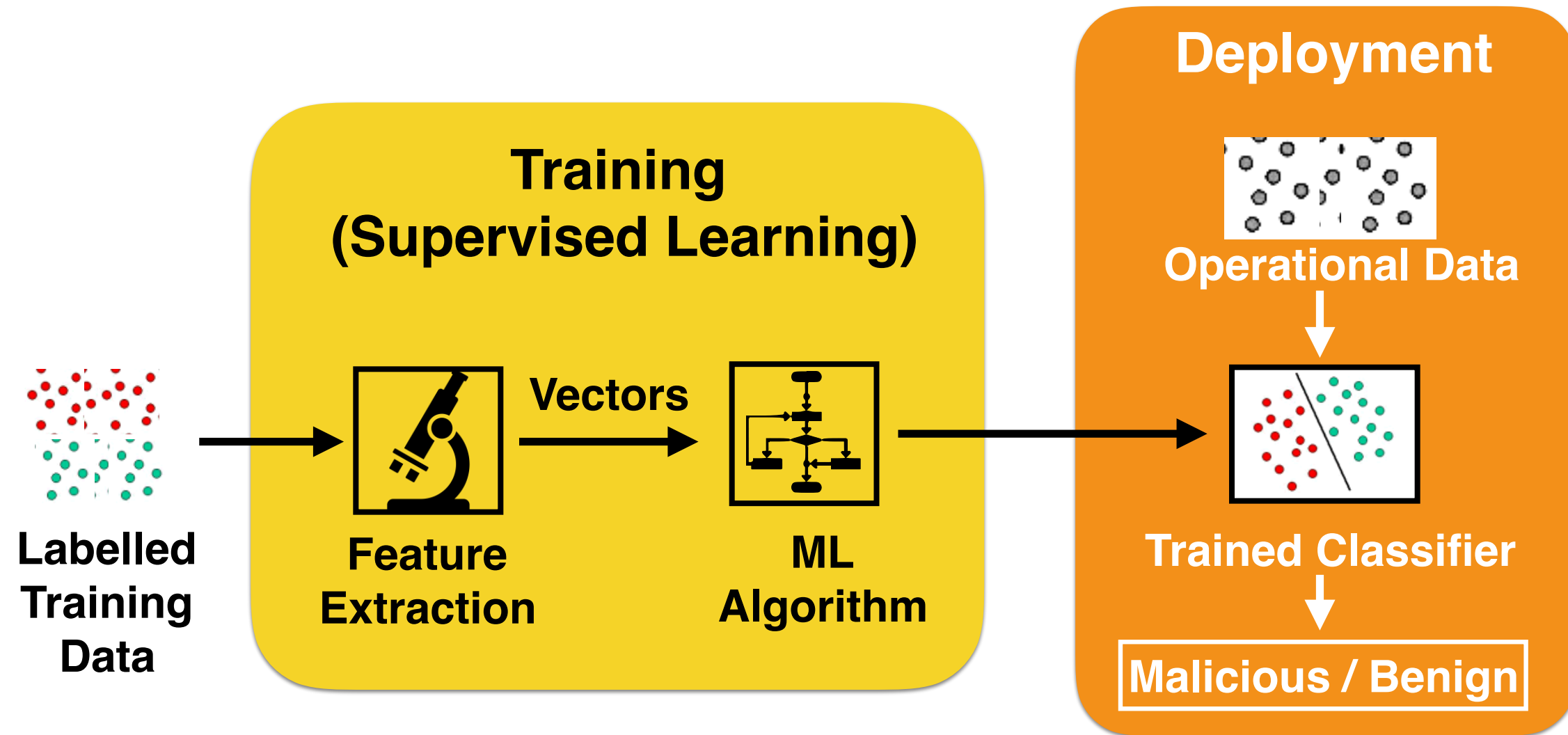# Security Tasks are Different: Adversary Adapts



**Goal**: Understand classifiers under attack.

**Results**: Vulnerable to automated evasion.

# Building Machine Learning Classifiers



**Labelled Training Data**

**Training (Supervised Learning)**

**Feature Extraction**

Vectors

**ML Algorithm**

**Trained Classifier**

# Assumption: Training Data is Representative



**Labelled Training Data**

**Training (Supervised Learning)**

Feature Extraction

Vectors

ML Algorithm

**Deployment**

Operational Data

Trained Classifier

**Malicious / Benign**

# Results: Evaded PDF Malware Classifiers

| | PDFrate* [ACSAC'12] | Hidost [NDSS'13] |
|---|---|---|
| Accuracy | 0.9976 | 0.9996 |
| False Negative Rate | 0.0000 | 0.0056 |
| False Negative Rate with Adversary | **1.0000** | **1.0000** |

\* Mimicus [Oakland '14], an open source reimplementation of PDFrate.

# Results: Evaded F清sifiers

> **Very robust against "strongest conceivable mimicry attack".**

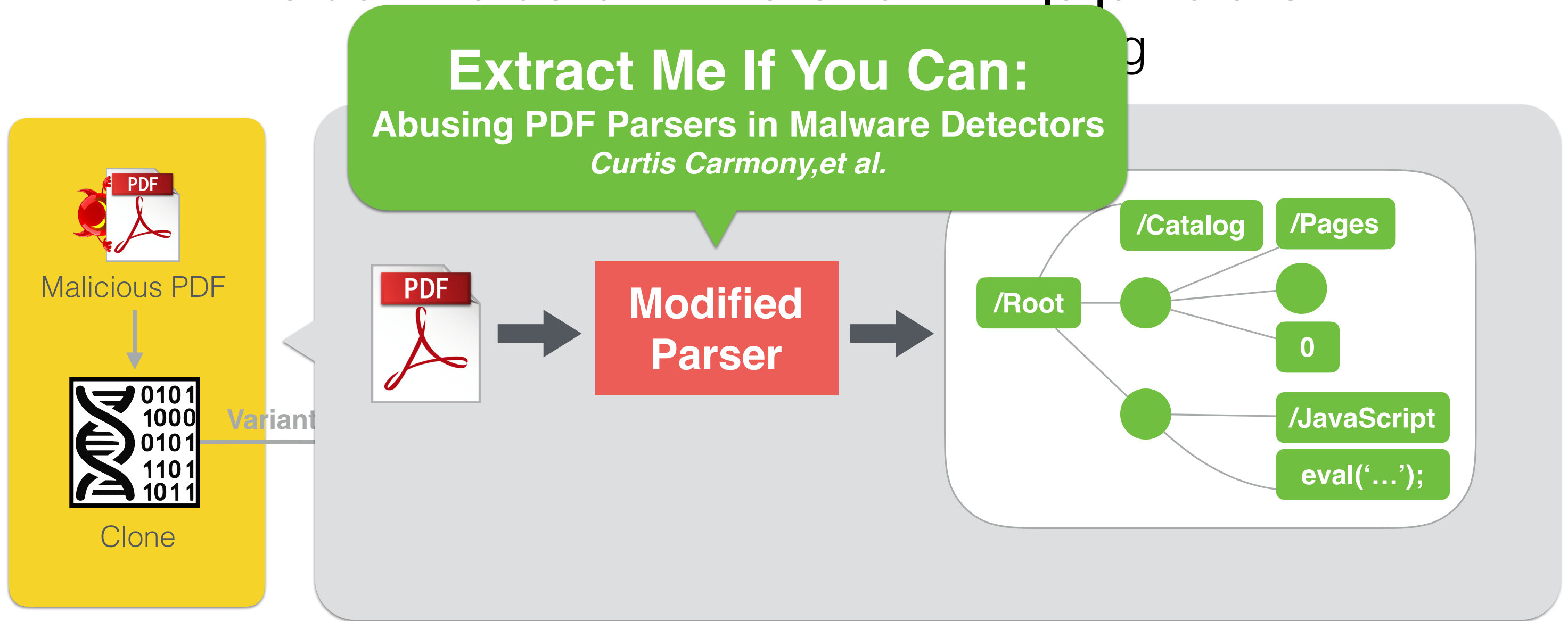| | PDFrate*<br>[ACSAC'12] | Hidost<br>[NDSS'13] |
|---|---|---|
| Accuracy | 0.9976 | 0.9996 |
| False Negative Rate | 0.0000 | 0.0056 |
| False Negative Rate with Adversary | **1.0000** | **1.0000** |

* Mimicus [Oakland '14], an open source reimplementation of PDFrate.

# Automated Evasion Approach
## Based on Genetic Programming



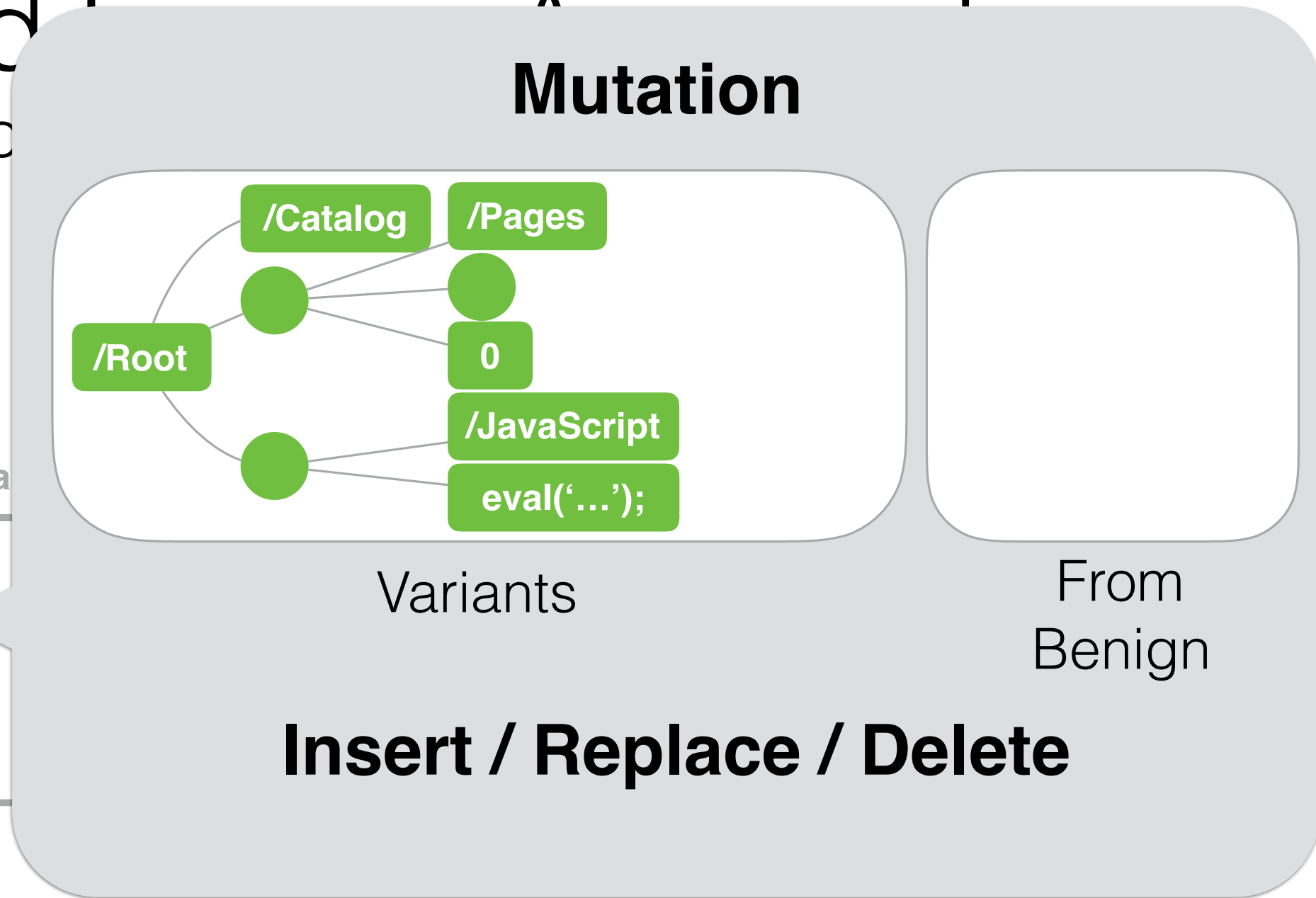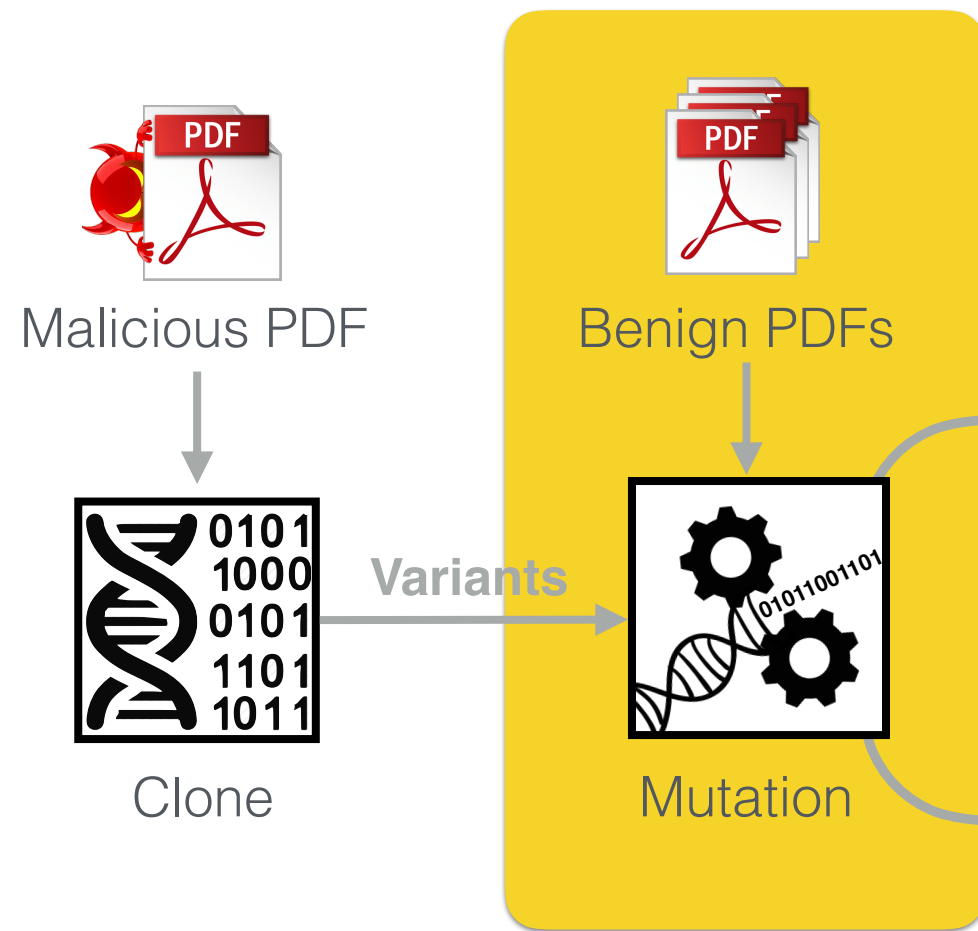Malicious PDF

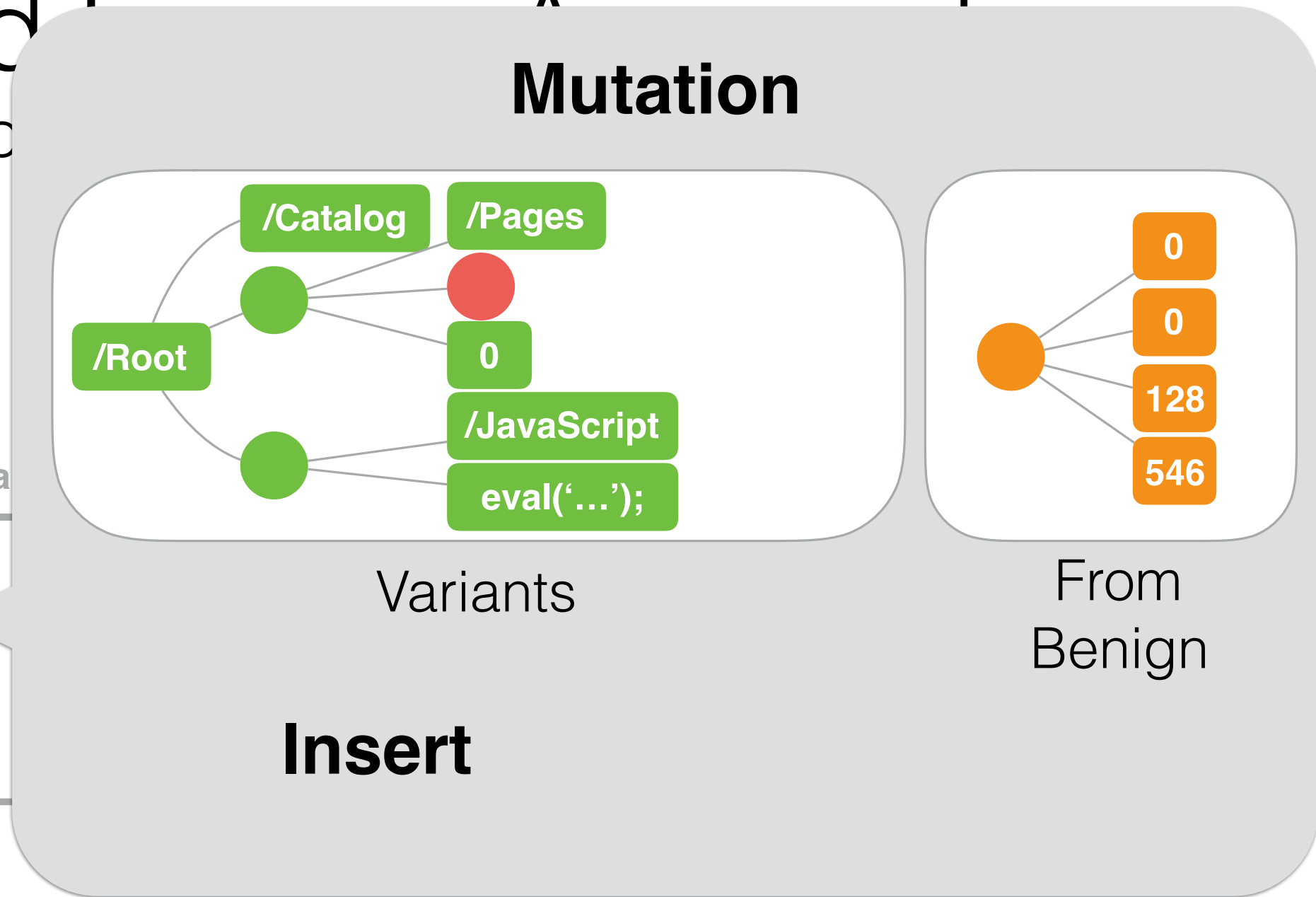Benign PDFs

**Variants**

Clone

**Variants**

Mutation

Select Variants

**Variants**

# Automated Evasion Approach



Malicious PDF

Clone

Variant

**Extract Me If You Can:**
**Abusing PDF Parsers in Malware Detectors**
*Curtis Carmony, et al.*

PDF → **Modified Parser** →

/Catalog  /Pages

/Root

0

/JavaScript

eval('...');

13

# Automated Fuzzing Approach

## Based on



Malicious PDF

Benign PDFs

Clone

Variants

Mutation

### Mutation

/Catalog /Pages

/Root 0

/JavaScript

eval('…');

Variants

From Benign

**Insert / Replace / Delete**

14

# Automated Evasion Attack

## Based on

**Malicious PDF**

**Benign PDFs**

Variants

**Variants**

Clone

Mutation

## Mutation

### Variants

/Catalog  /Pages

/Root

0

/JavaScript

eval('…');

### Insert

### From Benign

0

0

128

546

15

# Automated Fuzzing Approach
## Based on

**Malicious PDF**

**Benign PDFs**

**Clone**

**Mutation**

**Variants**

**Mutation**

/Catalog | /Pages

/Root

0

/JavaScript

eval('…');

0

0

128

546

**Variants**

**Insert**

**From Benign**

# Automated Evasion Approaches

## Based on



Malicious PDF

Clone

Benign PDFs

**Variants**

Mutation

### Mutation

/Catalog /Pages

/Root

0

0

0

128

546

/JavaScript

eval('...');

Variants

### Replace

From Benign

0

128

# Automated Fuzzing Approach

## Based on



Malicious PDF

Clone

Benign PDFs

Variants

Mutation

**Mutation**

/Catalog    /Pages

/Root

0

/JavaScript

eval('…');

0

128

Variants

**Replace**

From Benign

18

# Automated Evasion Approach

Based on



Malicious PDF

Benign PDFs

**Variants**

Clone

Mutation

## Mutation

/Catalog   /Pages

/Root   0

/JavaScript

eval('…');
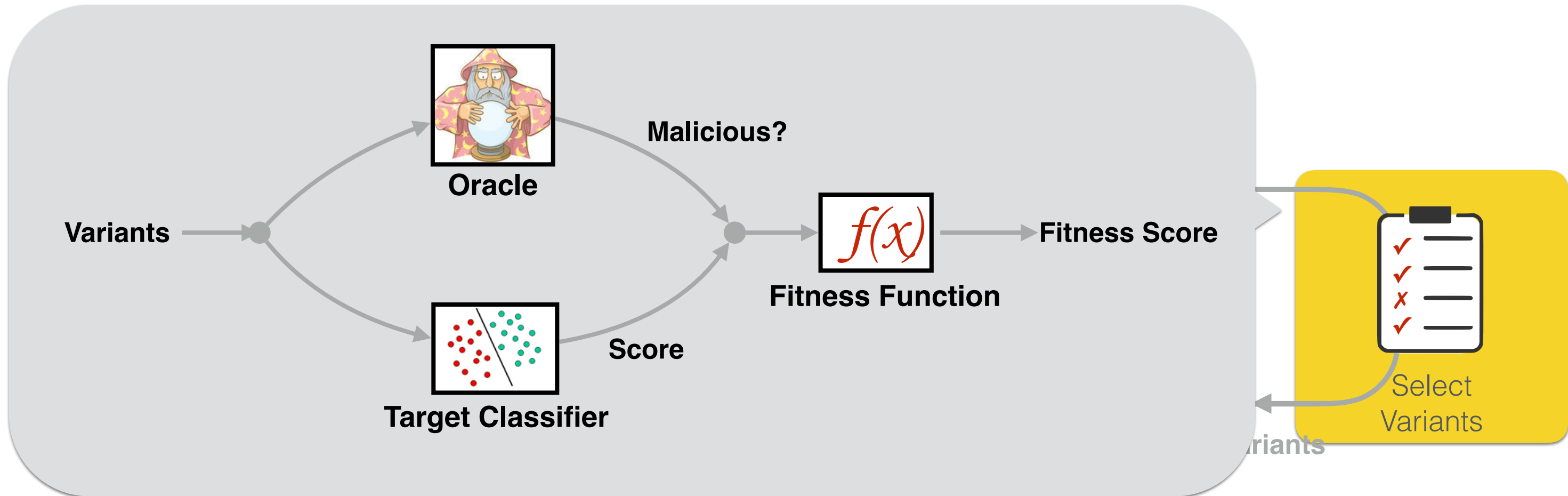
0
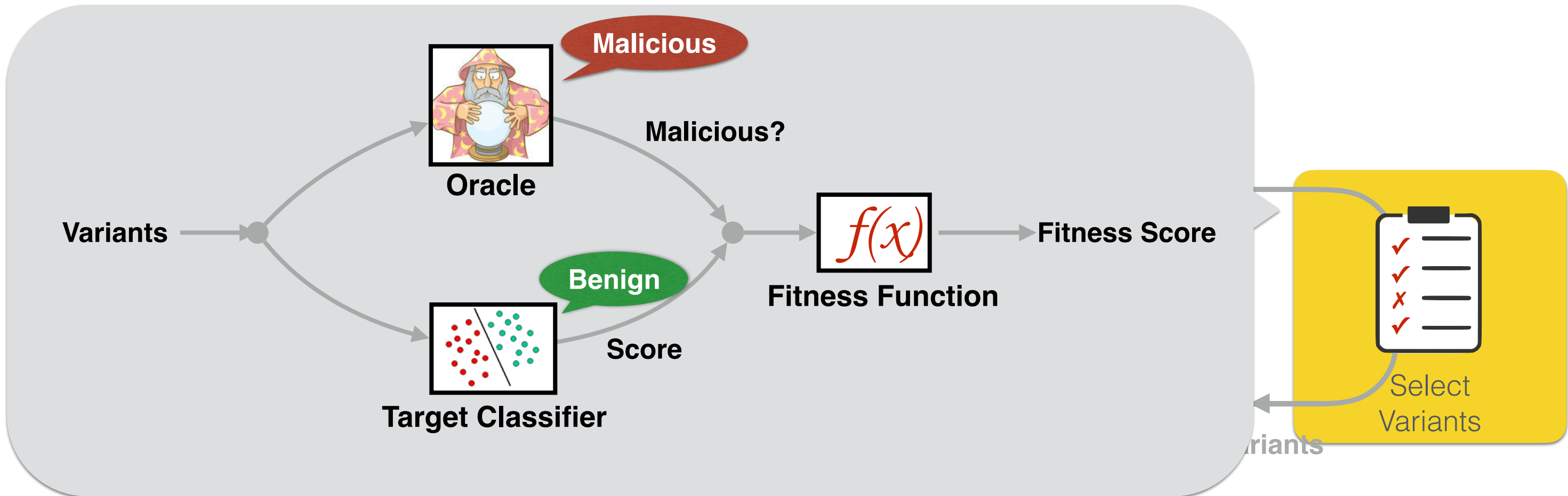
128

Variants

**Delete**

From Benign

19

# Automated Evasion Approach
## Based on Genetic Programming

# Automated Evasion Approach
## Based on Genetic Programming

# Automated Evasion Approach
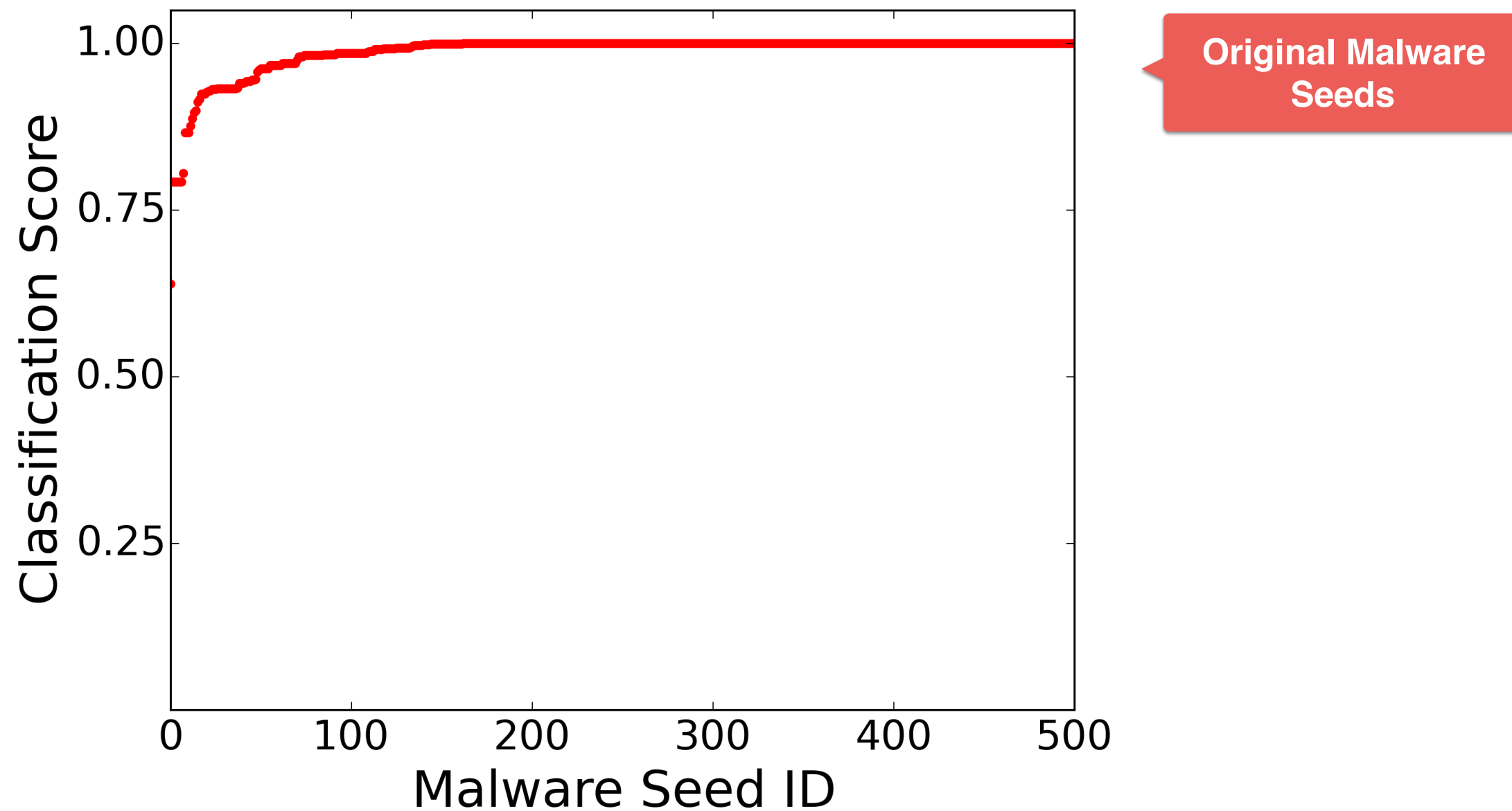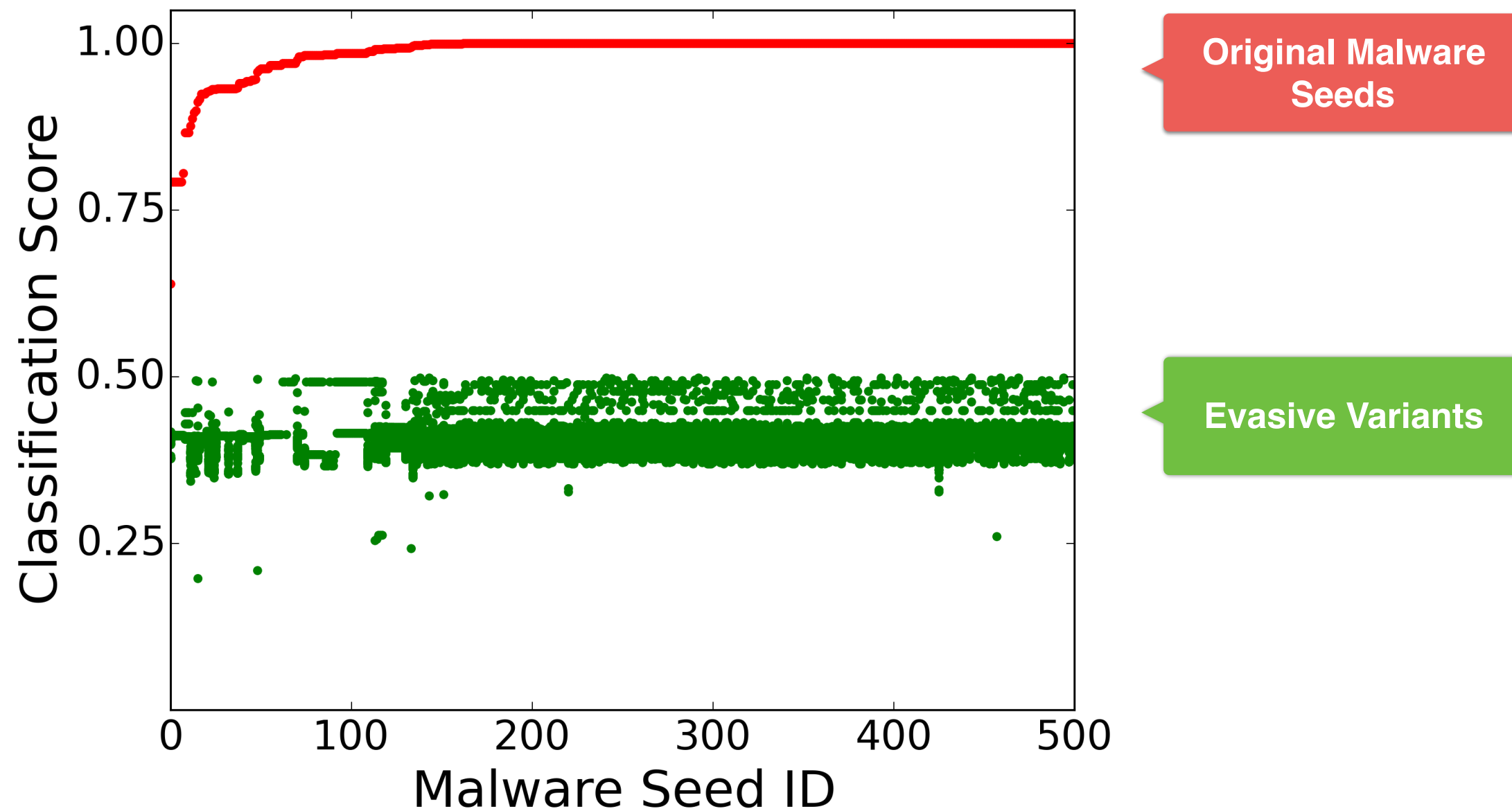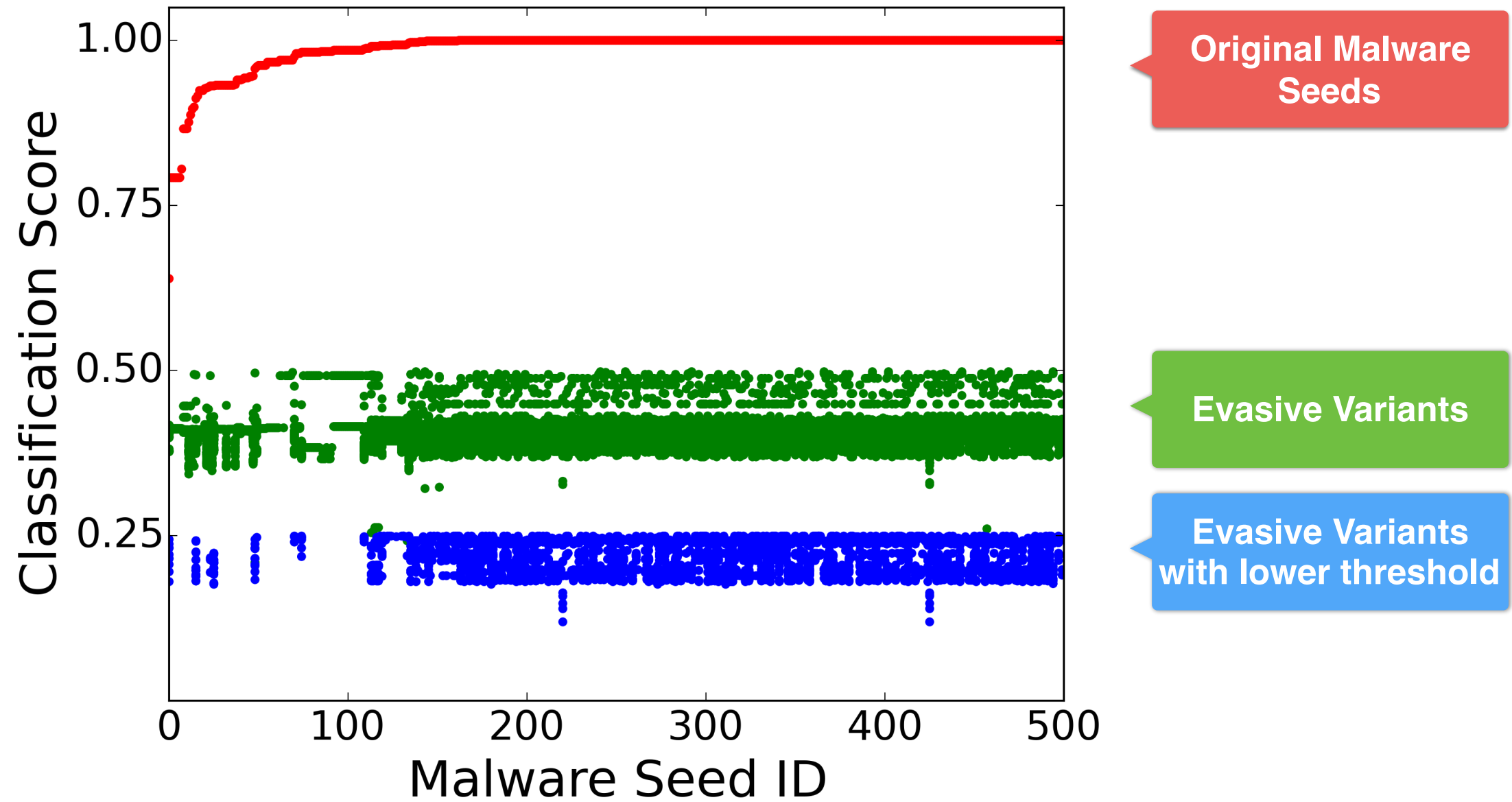## Based on Genetic Programming

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

Variants

Clone

Variants

Mutation

Variants

Select Variants

# Results: Evaded PDFrate 100%



Original Malware Seeds

# Results: Evaded PDFrate 100%

# Evaded PDFrate with Adjusted Threshold



**Original Malware Seeds**

**Evasive Variants**

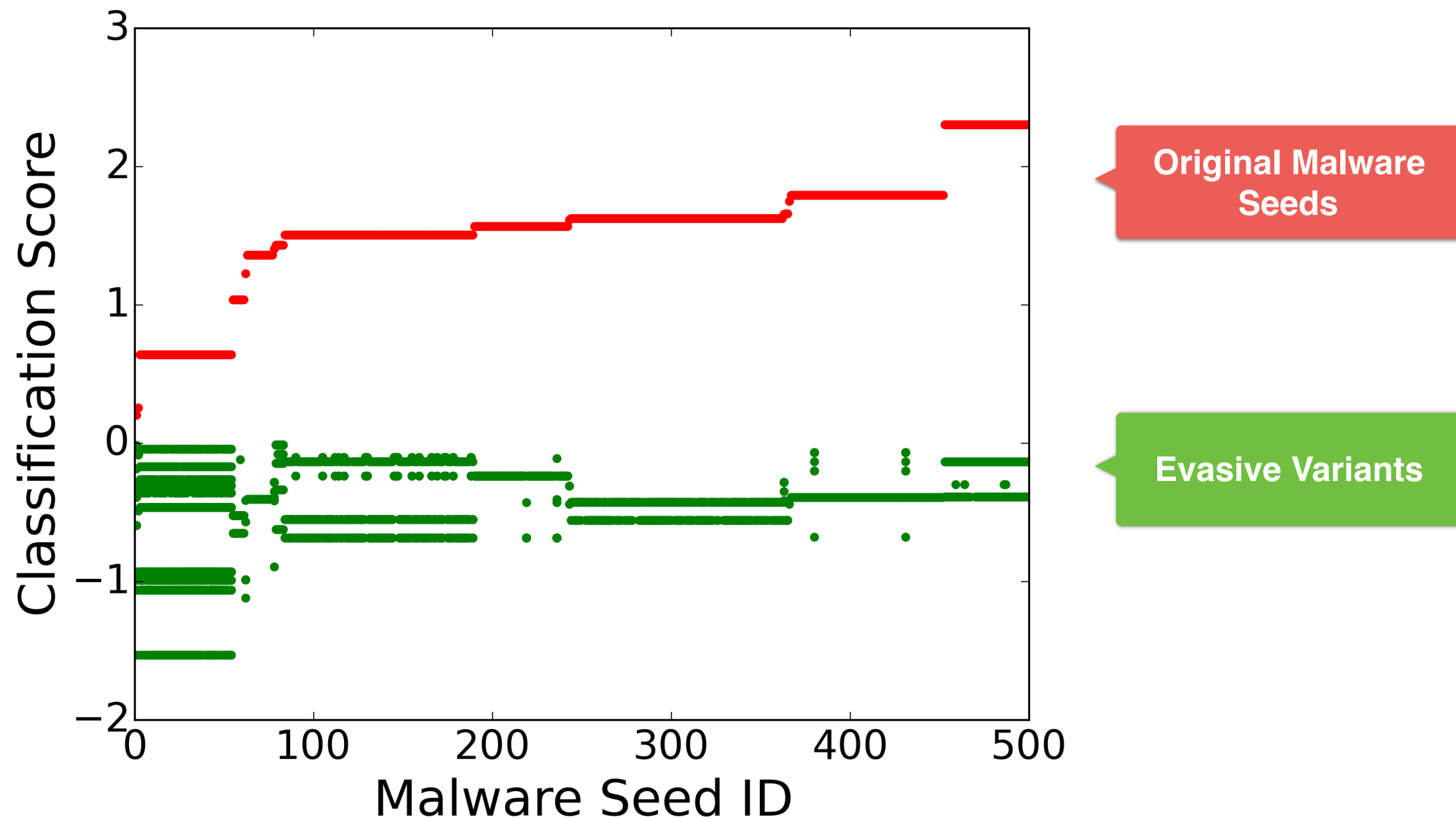**Evasive Variants with lower threshold**
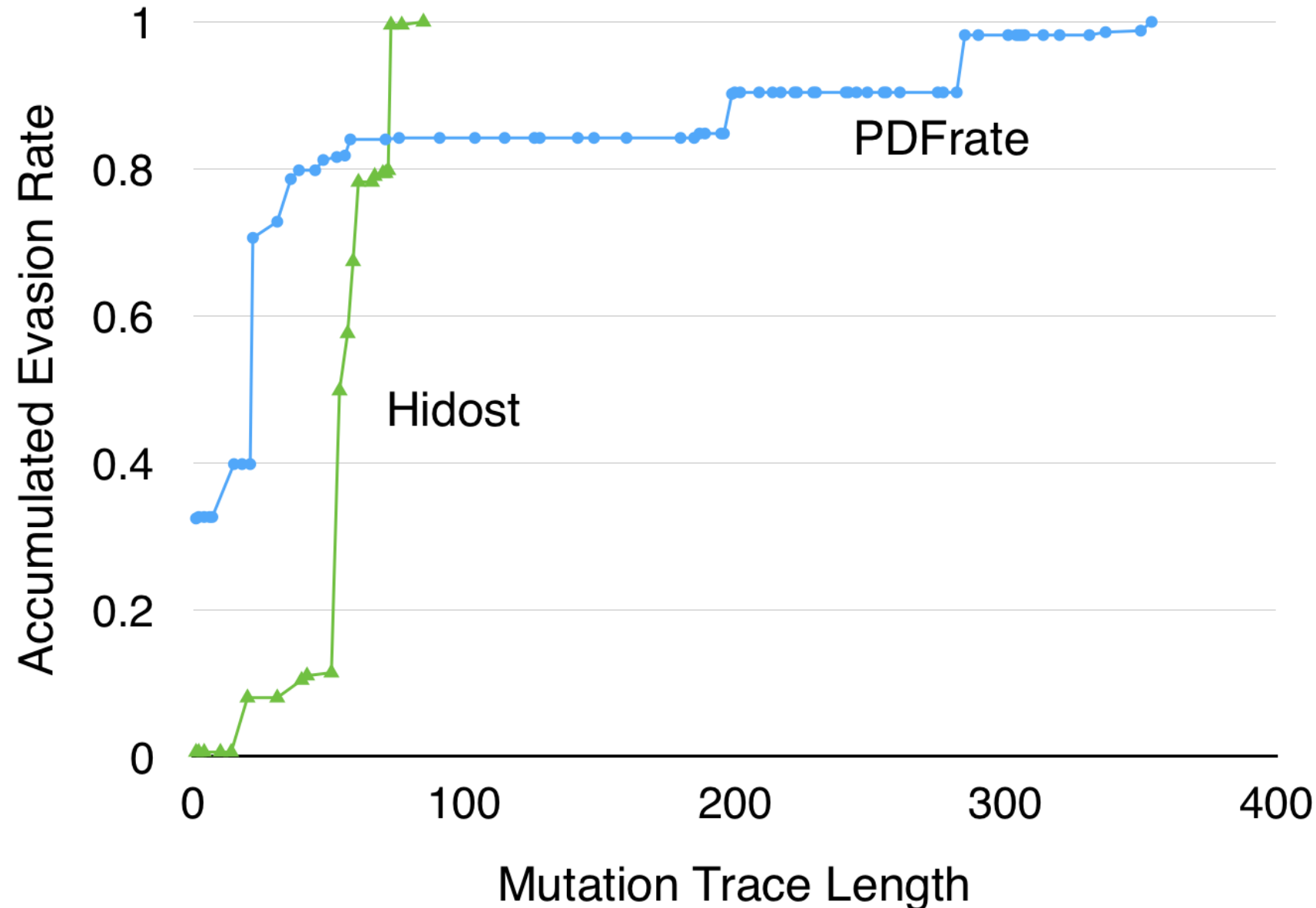
# Results: Evaded Hidost 100%



Original Malware Seeds

# Results: Evaded Hidost 100%

# Results: Accumulated Evasion Rate



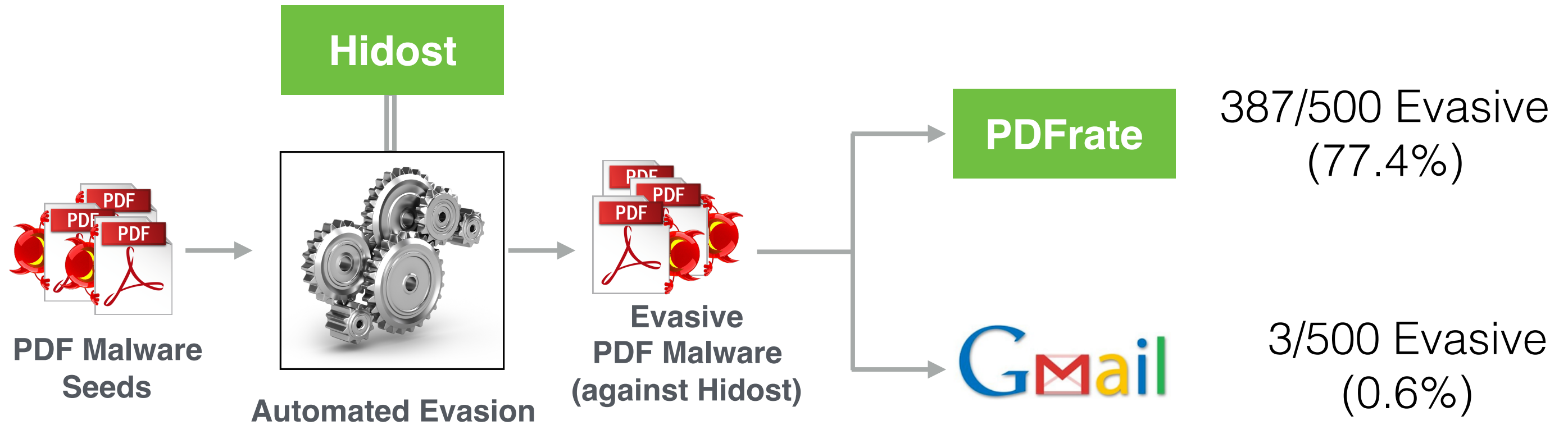**Difficulty varies by seed**
  Simple mutations often work
  Complex mutations sometimes
    needed.

**Difficulty varied by targets:**
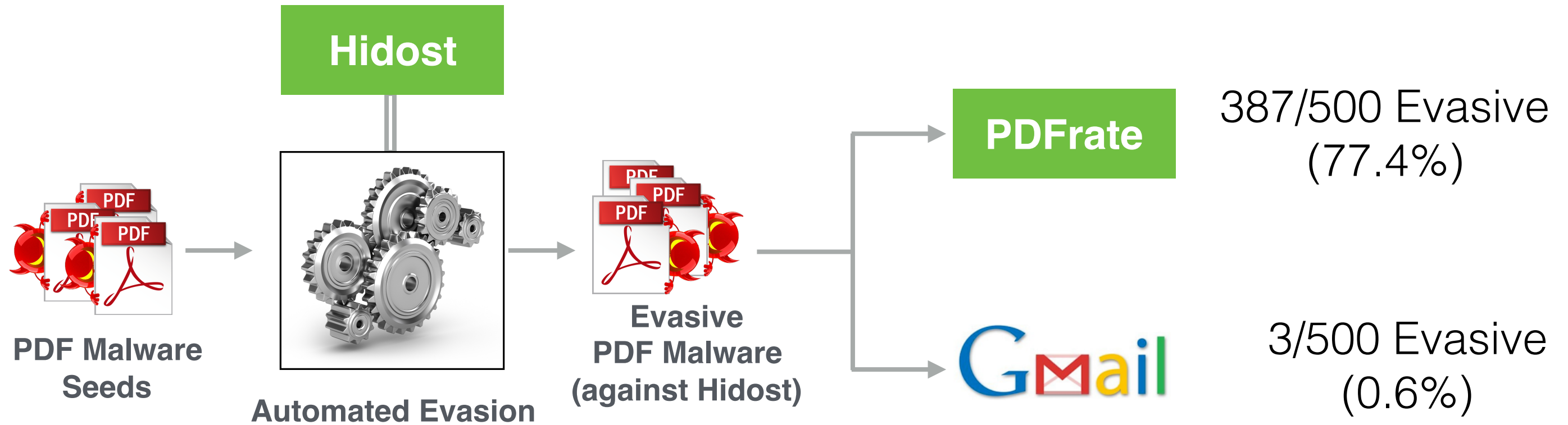  PDFrate: 6 days to evade all
  Hidost: 2 days to evade all

# Cross-Evasion Effects



**Hidost**

**PDFrate**

387/500 Evasive (77.4%)

3/500 Evasive (0.6%)

**PDF Malware Seeds**

**Automated Evasion**

**Evasive PDF Malware (against Hidost)**

**Gmail's classifier is secure?**

# Cross-Evasion Effects



**Hidost**

PDF Malware Seeds

**Automated Evasion**

Evasive
PDF Malware
(against Hidost)

**PDFrate**

387/500 Evasive
(77.4%)

Gmail

3/500 Evasive
(0.6%)

**Gmail's classifier is ~~secure?~~ different.**

# Evading Gmail's Classifier

```
1  for javascript in pdf.all_js:
2      javascript.append_code("var ndss=1;")
```
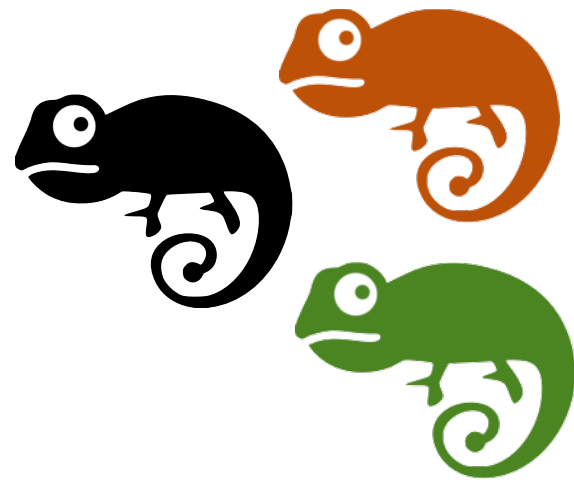
Evasion rate on Gmail : 135/380 (35.5%)
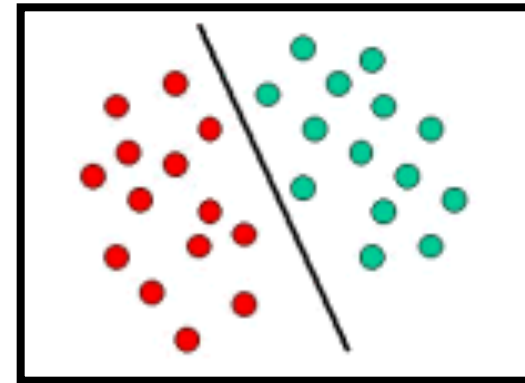
# Evading Gmail's Classifier

```python
1 for javascript in pdf.all_js:
2     javascript.append_code("var ndss=1;")
3
4 if pdf.get_size() < 7050000:
5    pdf.add_padding(7050000 - pdf.get_size())
```

Evasion rate on Gmail : 179/380 (47.1%)

# Conclusion



Vs.



# Who will win this arm race?

Source Code: http://EvadeML.org