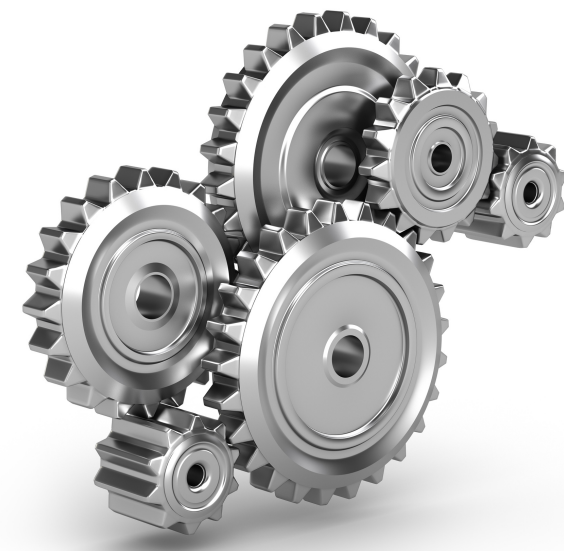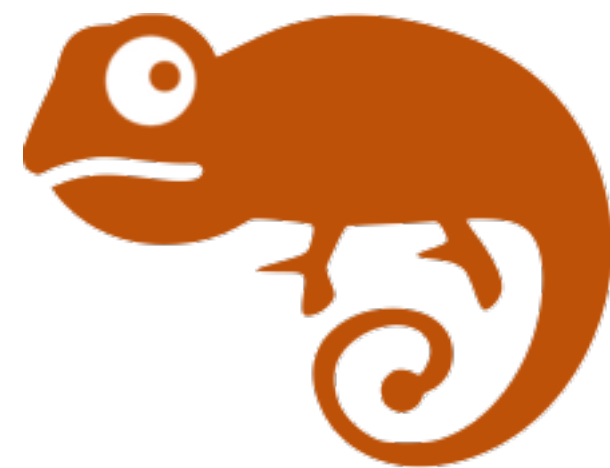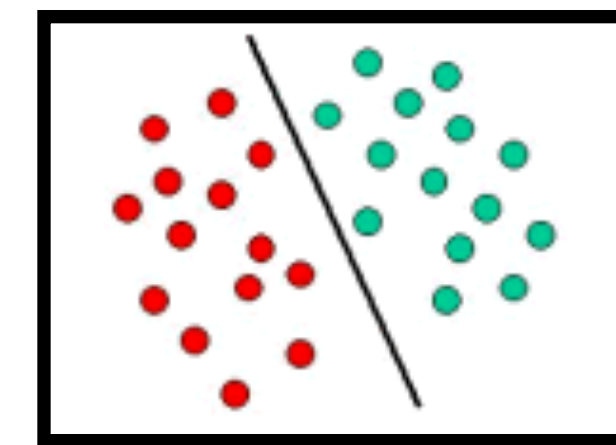# Automatically Evading Classifiers
## A Case Study on PDF Malware Classifiers

Weilin Xu

David Evans

Yanjun Qi

University of Virginia

# Machine Learning is Solving Our Problems

Spam

IDS

Fake Accounts

Malware

…

…

**Completed • $16,000 • 377 teams**

# 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) |

3

# 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

# Assumption: Training Data is Representative

# 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ssifiers

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

| | 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.

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

**Variants**

**Variants**

Clone

Mutation

Select Variants

**Variants**

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

Variants

Clone

**Variants**

Mutation
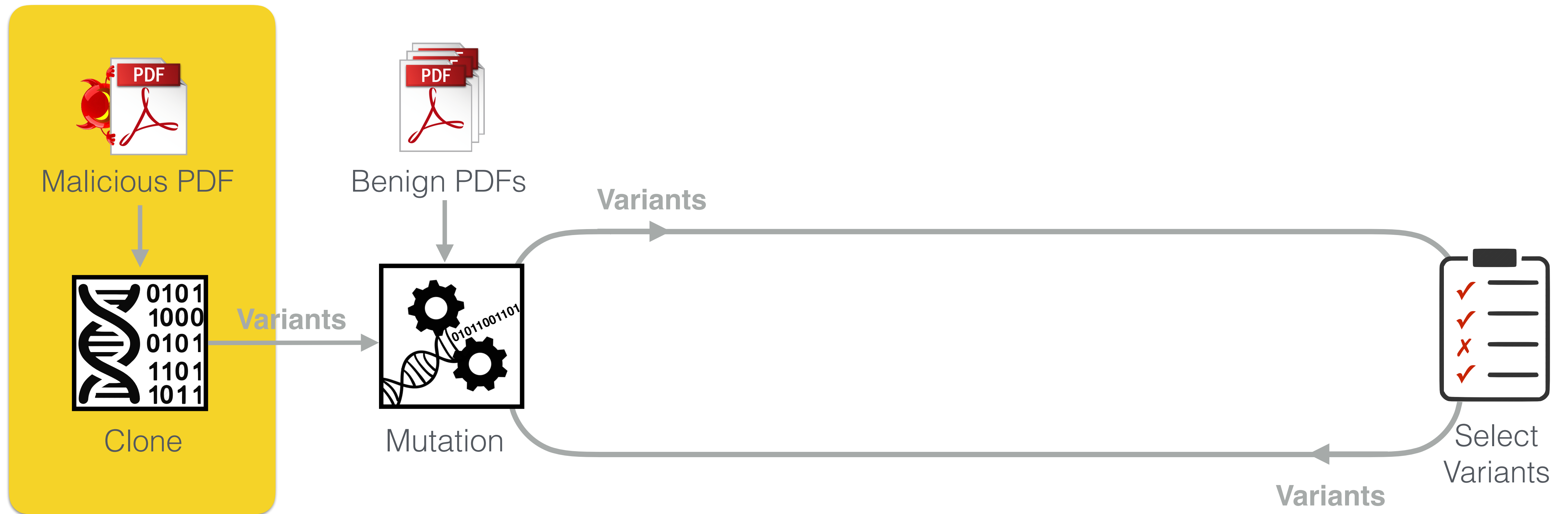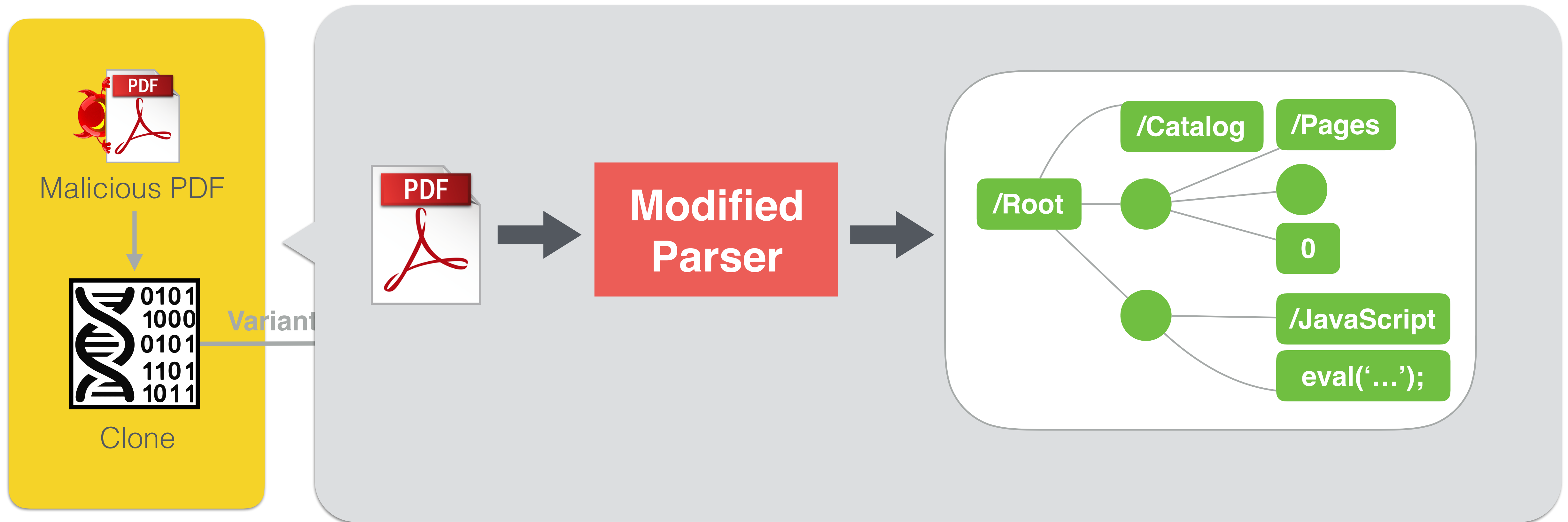
**Variants**

Select Variants

# Automated Evasion Approach
## Based on Genetic Programming

# Automated Evasion Approach
## Based on Genetic Programming

# Automated Evasion Approach



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

Malicious PDF

Variant

Clone

PDF

**Modified Parser**

/Catalog   /Pages

/Root

0

/JavaScript

eval('…');

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Clone

Benign PDFs

**Variants**

Mutation

**Variants**

**Variants**

Select Variants

**Variants**

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

Clone

**Variants**

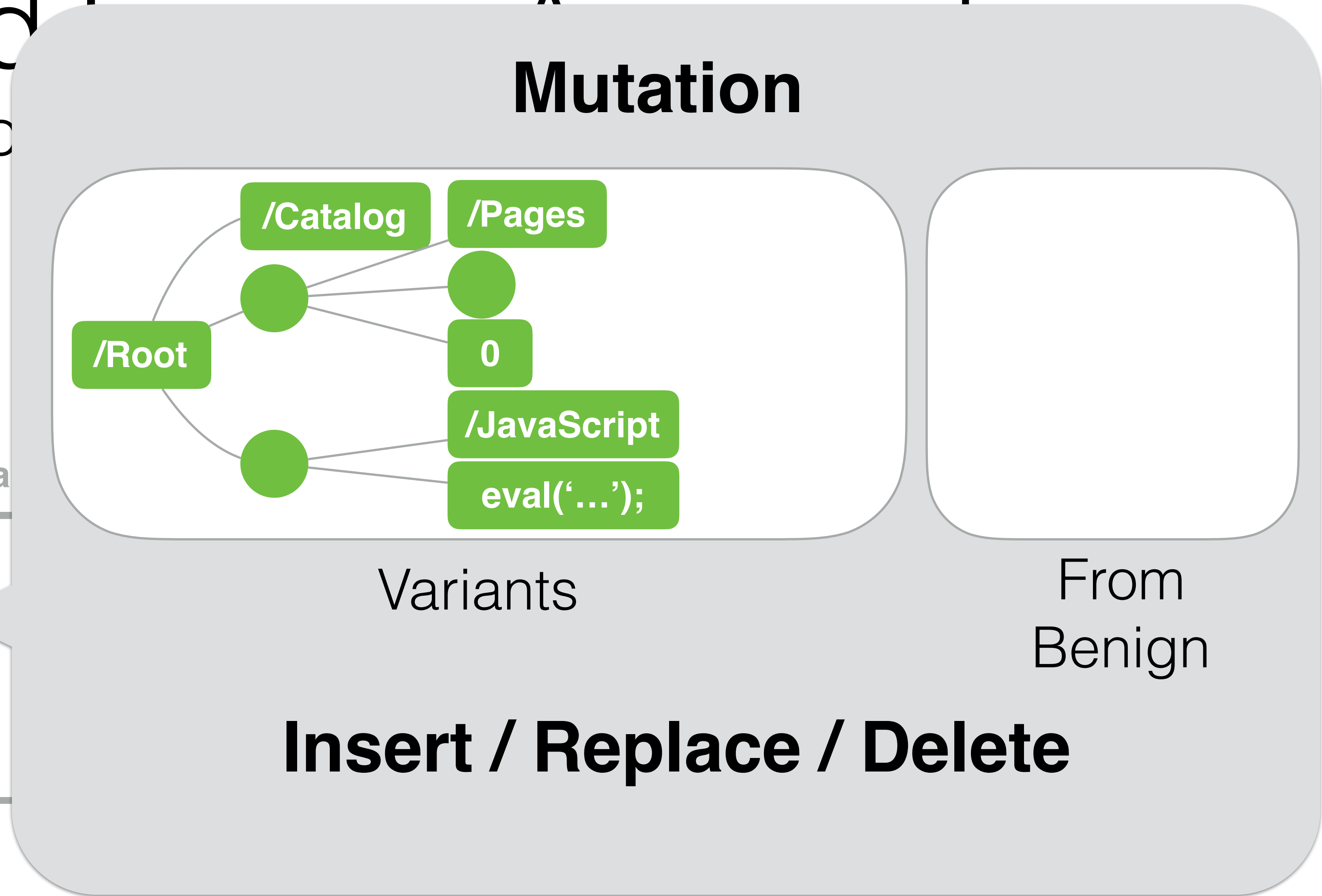Mutation

**Variants**

**Variants**

Select Variants

# Automated Fuzzing Approach

## Based on

Malicious PDF

Benign PDFs

**Mutation**

/Catalog   /Pages

/Root   0

/JavaScript

eval('...');

Variants

From Benign

**Insert / Replace / Delete**

Variants

Clone   Mutation

# Automated Engine Approach

## Based on



Malicious PDF

Benign PDFs

Clone

Mutation

**Variants**

## Mutation

| /Catalog | /Pages |

/Root

0

/JavaScript

eval('…');

Variants

From Benign

**Insert**

# Automated Evasion Approach

## Based o...



Malicious PDF

Benign PDFs

Clone

Mutation

**Mutation**

/Catalog /Pages

/Root 0

/JavaScript

eval('…');

Variants

**Insert**

From Benign

21

# Automated Fuzzing Approach

## Based on



Malicious PDF

Benign PDFs

**Variants**

Clone

Mutation

### Mutation

**Mutation**

/Catalog  /Pages

/Root

0

/JavaScript

eval('…');

Variants

**Insert**

0

0

128

546

From Benign

22

# Automated Evasion Attack

## Based on



Malicious PDF

Benign PDFs

Clone

Mutation

Variants

**Mutation**

/Catalog    /Pages    0

/Root    0    0

/JavaScript    128

eval('…');    546

Variants

**Insert**

From Benign

23

# Automated Fuzzing Approach

## Based on



Malicious PDF

Benign PDFs

**Variants**

Clone

Mutation

### Mutation

/Catalog /Pages

/Root

0

/JavaScript

eval('…');

0

0

128

546

Variants

Insert

From Benign

24

# Automated Engine Approach

## Based on

Malicious PDF

Benign PDFs

Clone

Variants

Mutation

### Mutation

/Catalog /Pages 0

/Root 0 0

0 128

/JavaScript 546

eval('...');

Variants

### Replace

From Benign

25

# Automated Evasion Approach

## Based o



Malicious PDF

Clone

Benign PDFs

Mutation

Variants

**Mutation**

/Catalog    /Pages

/Root

0

/JavaScript

eval('…');

0
0
128
546

Variants

**Replace**

From
Benign

Automated Fuzzing Approach

Based on

Malicious PDF

Clone

Variants

Benign PDFs

Mutation

**Mutation**

/Catalog    /Pages    0    0    0    128    546

/Root

0

/JavaScript

eval('…');

Variants

**Replace**

From
Benign

# Automated Evasion Approach

## Based on

**Malicious PDF**

**Benign PDFs**

**Variants**

**Clone**

**Mutation**

### Mutation

/Catalog   /Pages

/Root

0

/JavaScript

eval('…');

0

128

**Variants**

**Replace**

**From Benign**

28

# Automated Evasion Approach

## Based on



Malicious PDF

Clone

**Variants**

Benign PDFs

Mutation

Variants

### Mutation

/Catalog  /Pages

/Root

0

/JavaScript

eval('…');

0

128

Variants

### Replace

From Benign

# Automated Engine: Approach

## Based on

Malicious PDF

Benign PDFs

**Variants**

Clone

Mutation



**Mutation**

/Catalog /Pages

/Root 0 0

128

/JavaScript

eval('…');

Variants

**Delete**

From Benign

30

# Automated Evasion Approach

## Based o

**Malicious PDF**

**Benign PDFs**

Clone

**Variants**

Mutation

### Mutation

/Catalog   /Pages

/Root

0

/JavaScript

eval('…');

0

128

Variants

### Delete

From Benign

31

Automated Fuzzing Approach

Based on

**Mutation**

/Catalog /Pages

/Root 0 0 128

eval('…');

Variants

Malicious PDF

Benign PDFs

Variants

Clone

Mutation

From Benign

**Delete**

# Automated Evasion Approach

## Based on

Malicious PDF

Benign PDFs

Clone

**Variants**

Mutation

### Mutation

/Catalog  /Pages

/Root

0

128

eval('…');

Variants

From Benign

**Insert / Replace / Delete**

33

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

**Variants**

Clone

**Variants**

Mutation

Select Variants

**Variants**

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

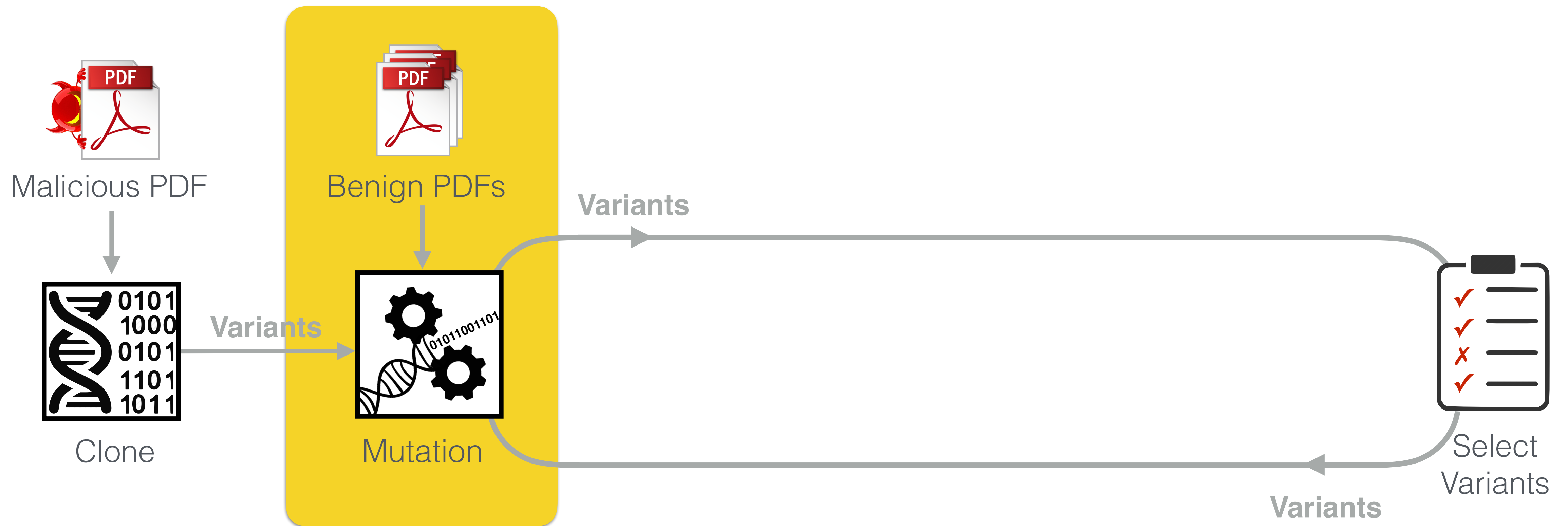Benign PDFs
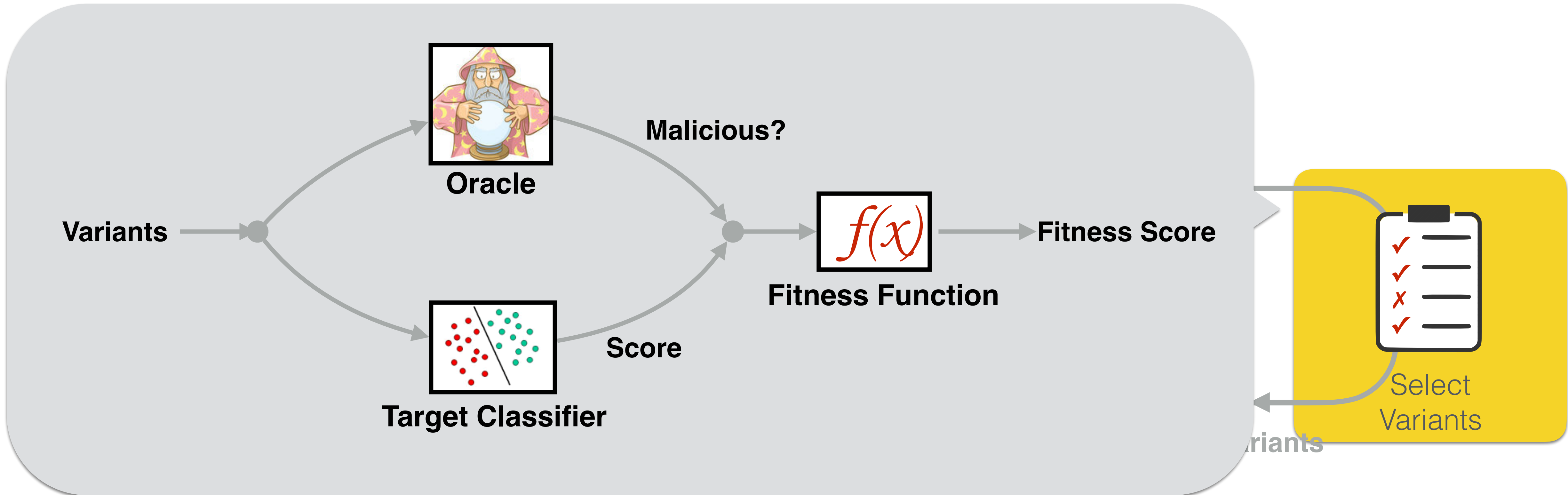
**Variants**

Clone

**Variants**

Mutation

**Variants**

Select Variants

# Automated Evasion Approach
## Based on Genetic Programming

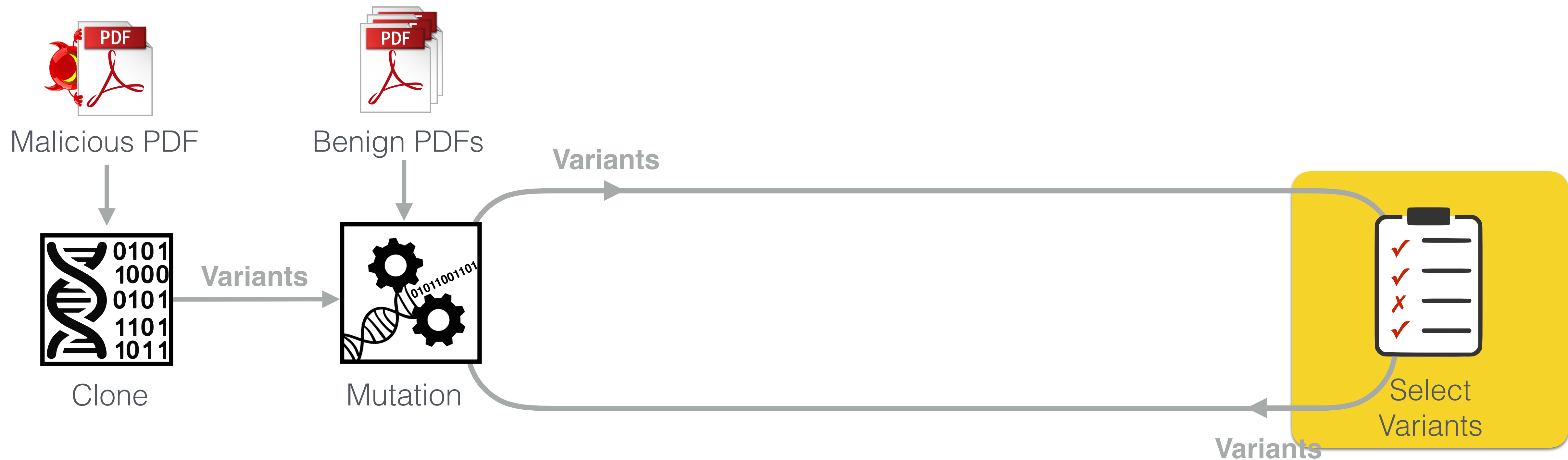# Automated Evasion Approach
## Based on Genetic Programming



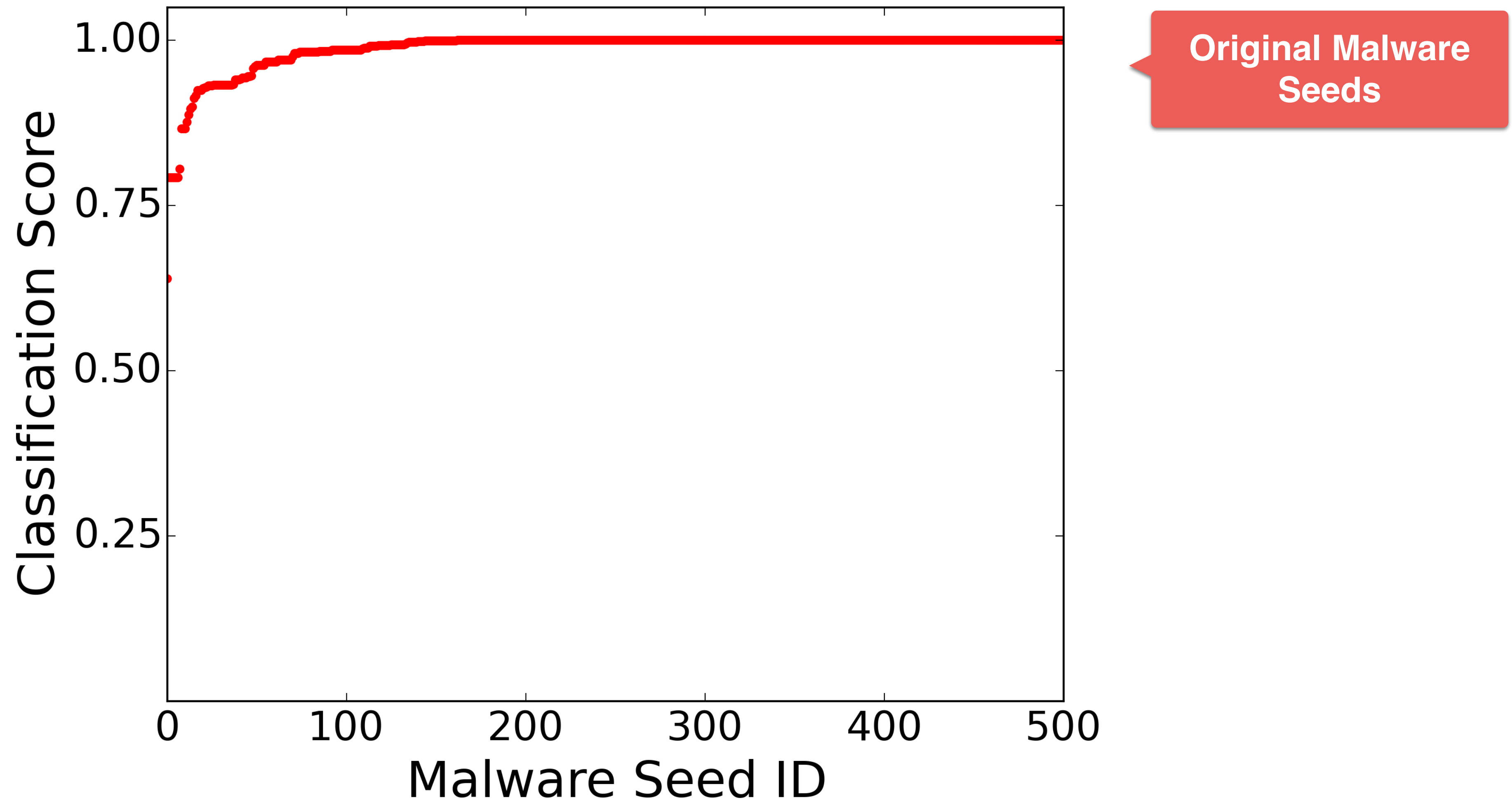Malicious PDF

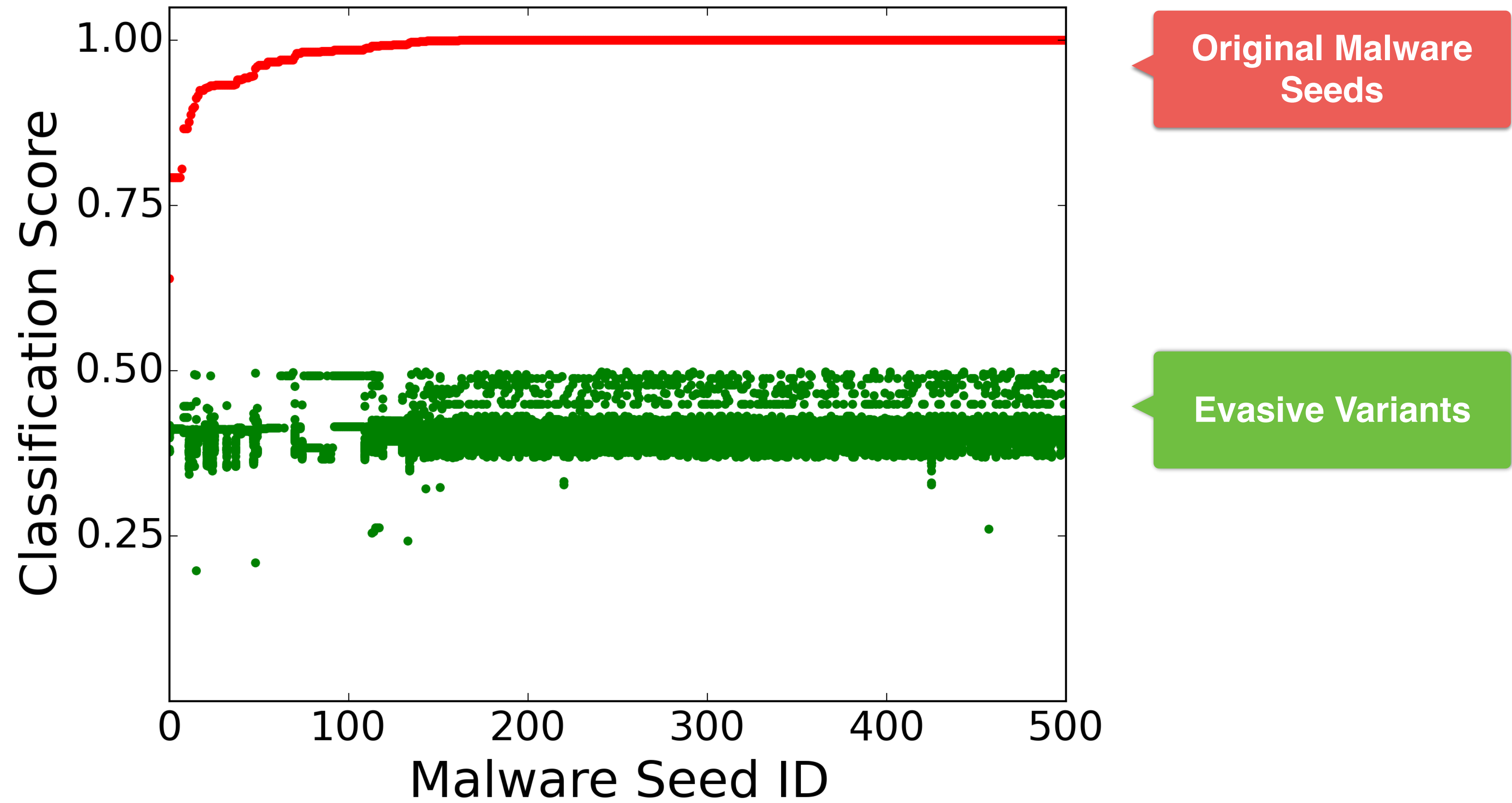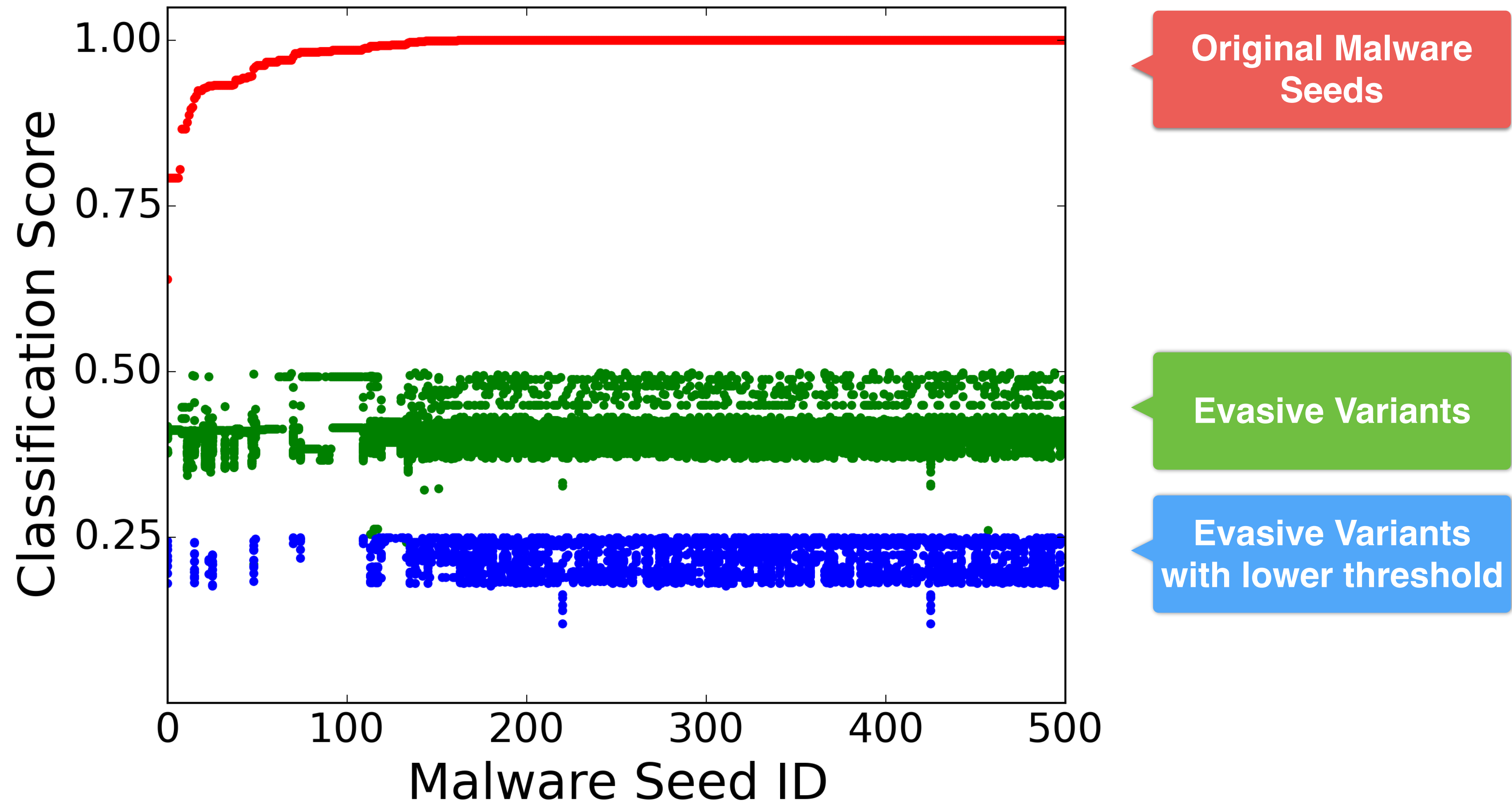Benign PDFs

Clone

**Variants**

Mutation

**Variants**

**Variants**

Select Variants

# Automated Evasion Approach
## Based on Genetic Programming



Malicious PDF

Benign PDFs

Variants

Variants

Clone

Mutation

Variants

Select
Variants

# Results: Evaded PDFrate 100%



Original Malware Seeds

# Evaded PDFrate with Adjusted Threshold



Original Malware Seeds

Evasive Variants

# Evaded PDFrate with Adjusted Threshold



Original Malware Seeds

Evasive Variants

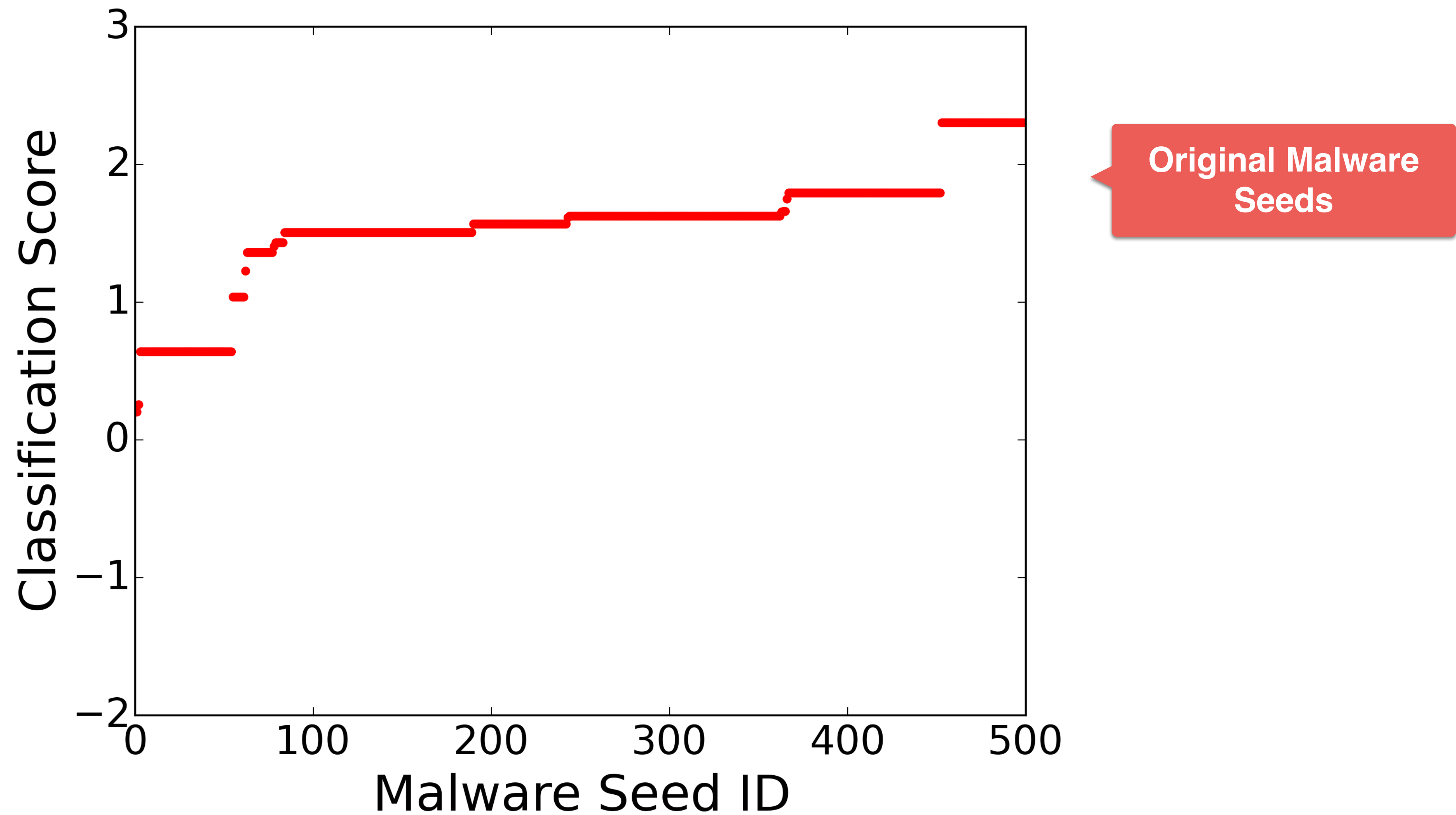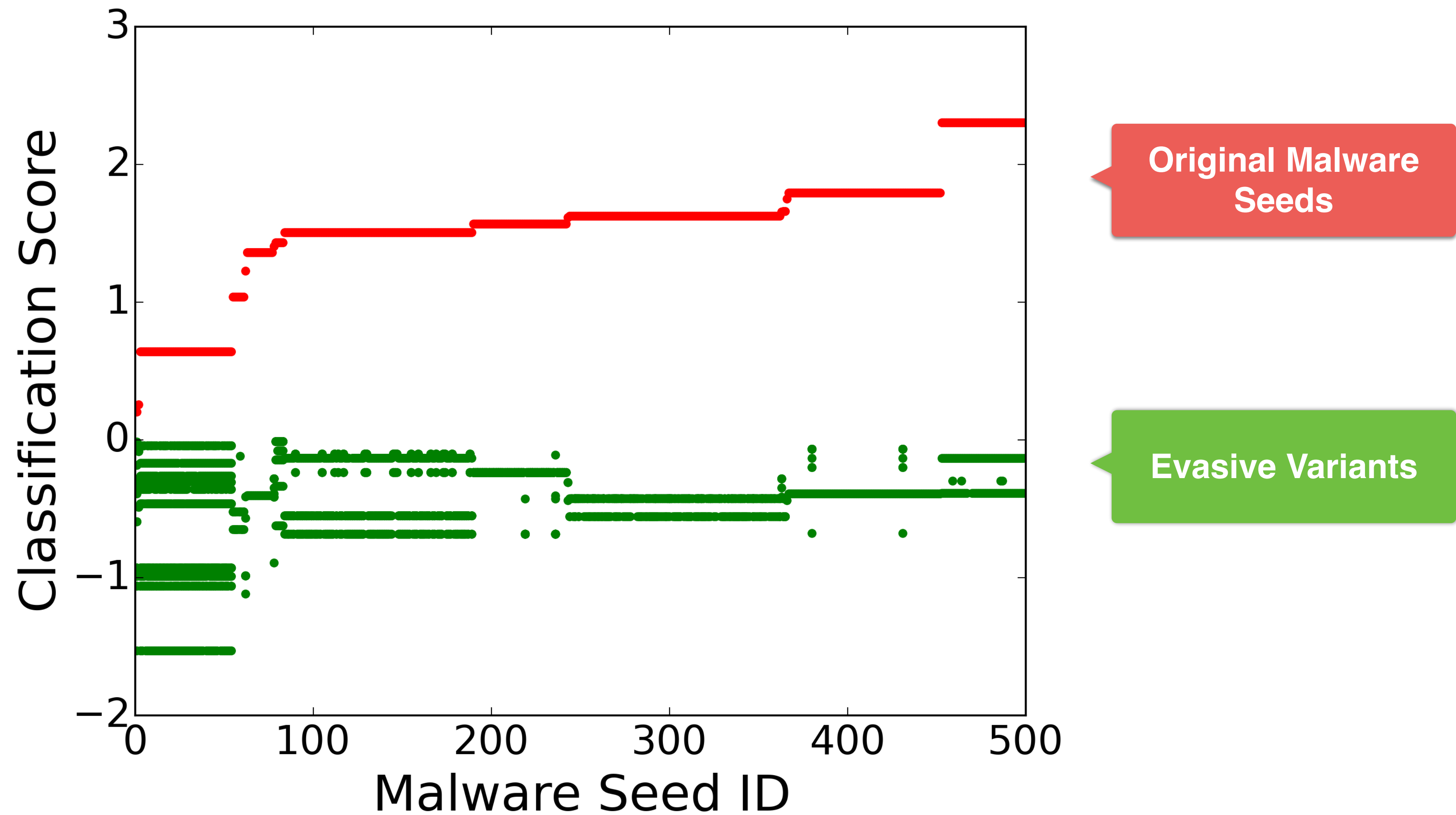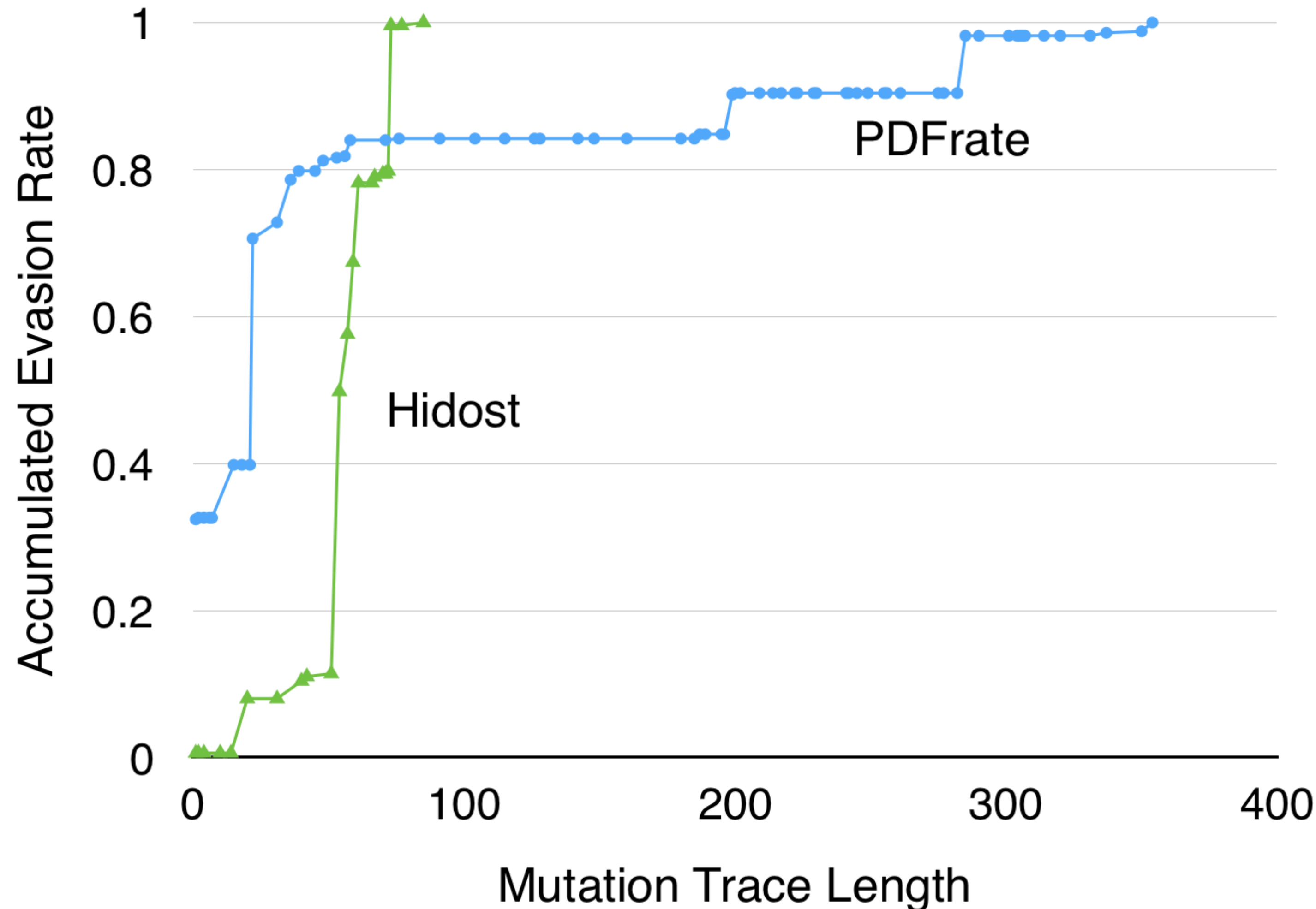Evasive Variants with lower threshold

# Results: Evaded Hidost 100%

# Results: Evaded Hidost 100%

# Results: Accumulated Evasion Rate



**Difficulties varied on targets.**
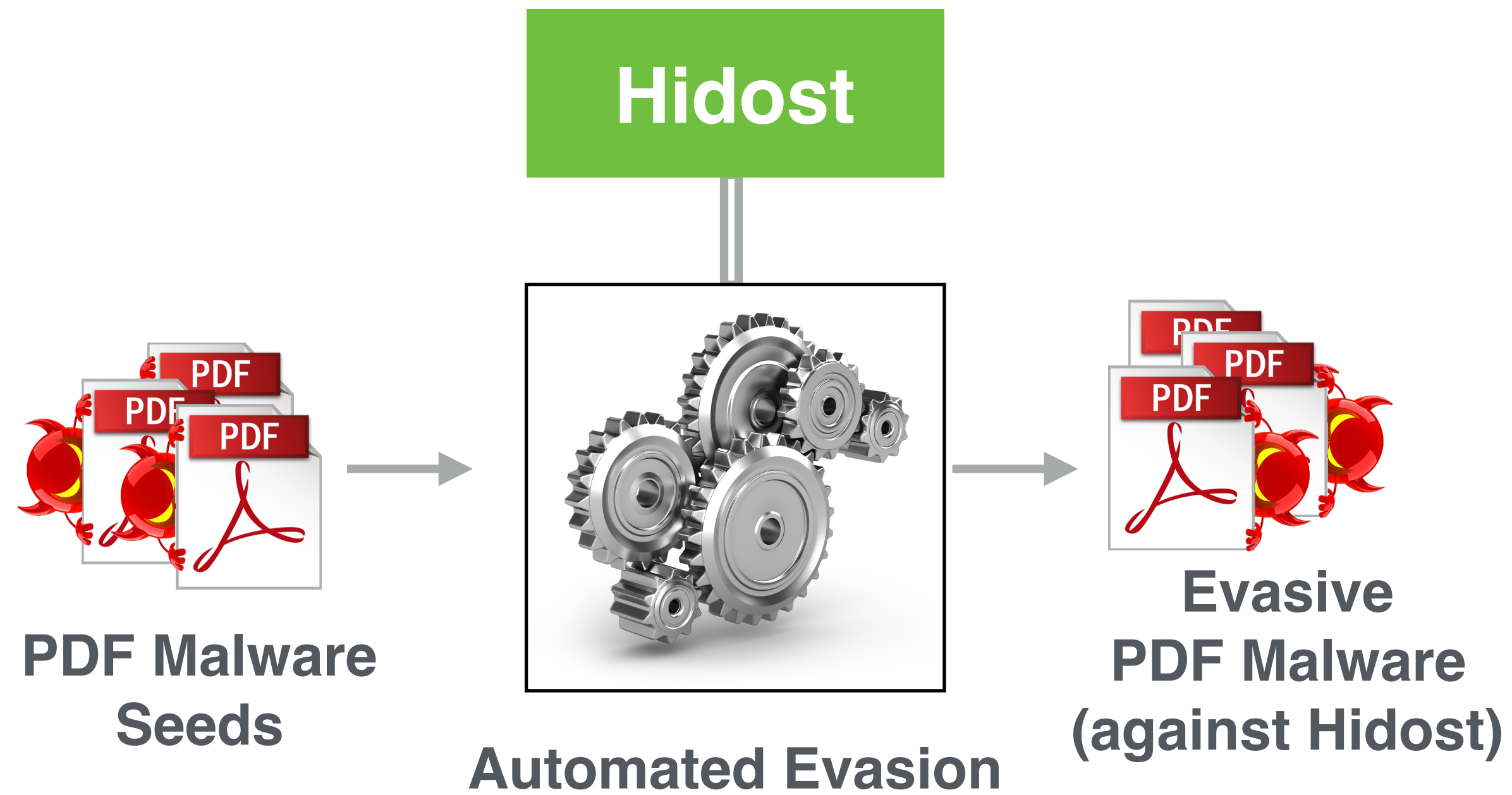    Evaded PDFrate in 6 days.
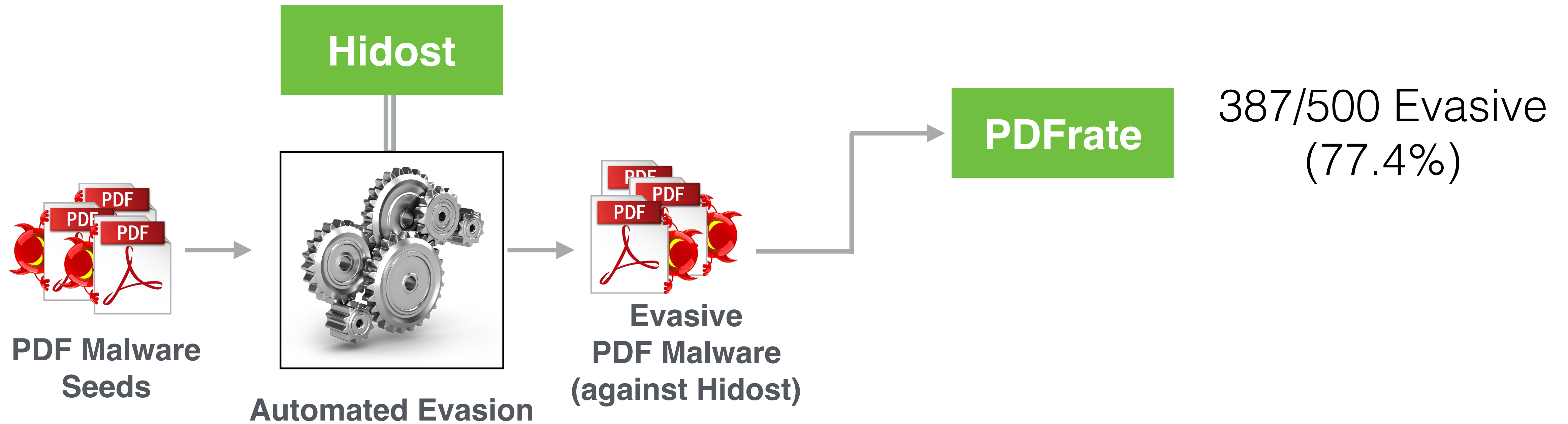    Evaded Hidost in 2 days.


**Difficulties varied on seeds.**
    Simple mutations worked.
    Complex mutations required.

# Cross-Evasion Effects



**PDF Malware Seeds**
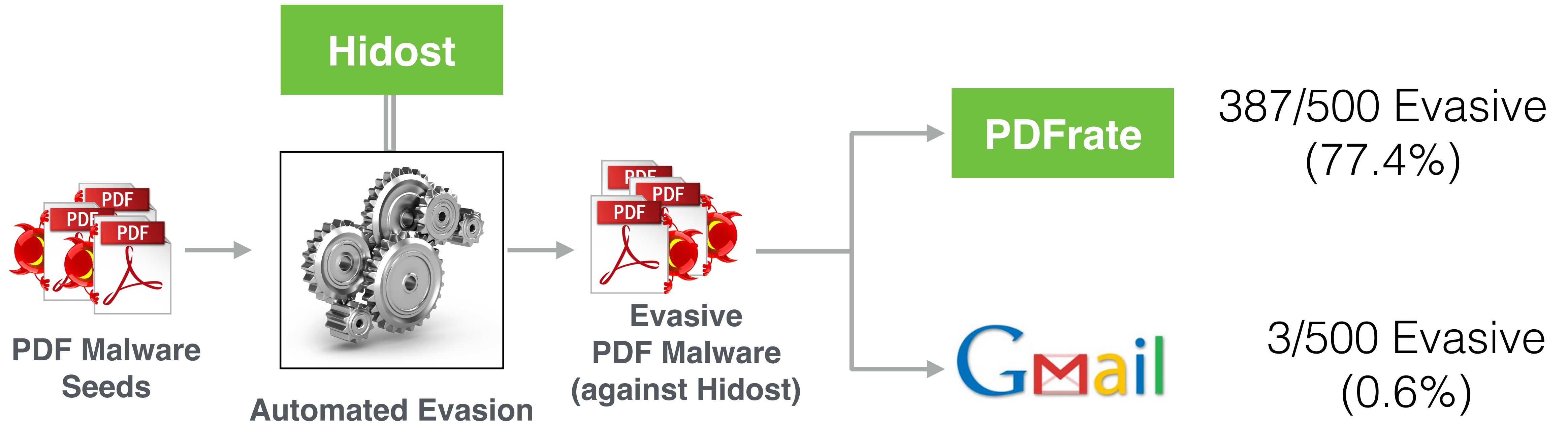
**Automated Evasion**

**Hidost**

**Evasive PDF Malware (against Hidost)**

# Cross-Evasion Effects



**Hidost**

**PDFrate**

387/500 Evasive
(77.4%)

**PDF Malware
Seeds**

**Automated Evasion**

**Evasive
PDF Malware
(against Hidost)**

# Cross-Evasion Effects



**PDF Malware Seeds** → **Hidost** / **Automated Evasion** → **Evasive PDF Malware (against Hidost)** → **PDFrate** / **Gmail**

387/500 Evasive (77.4%)

# Cross-Evasion Effects

**Hidost**

**PDFrate**

387/500 Evasive
(77.4%)

**PDF Malware
Seeds**

**Automated Evasion**

**Evasive
PDF Malware
(against Hidost)**

3/500 Evasive
(0.6%)

# Cross-Evasion Effects



**Hidost**

**PDFrate**

387/500 Evasive
(77.4%)

**PDF Malware
Seeds**

**Automated Evasion**

**Evasive
PDF Malware
(against Hidost)**

3/500 Evasive
(0.6%)

**Gmail's classifier is secure?**

# Cross-Evasion Effects

**Hidost**

**PDFrate**

387/500 Evasive
(77.4%)

**PDF Malware
Seeds**

**Automated Evasion**

**Evasive
PDF Malware
(against Hidost)**

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;")
```
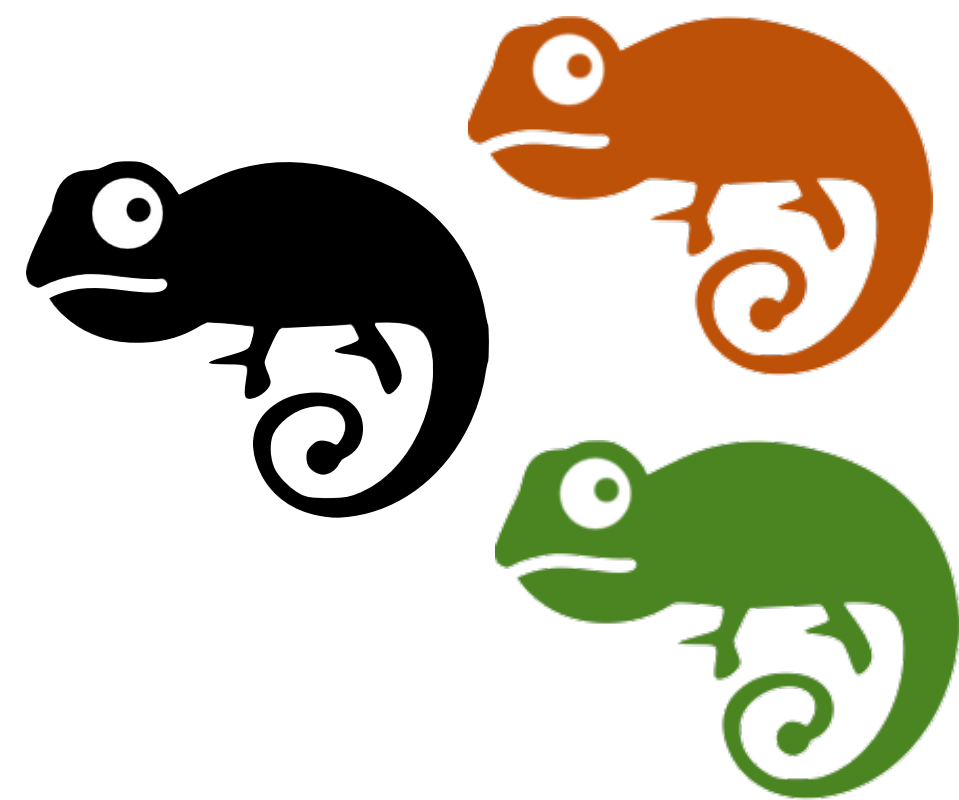
Evasion rate on Gmail : 135/380 (35.5%)

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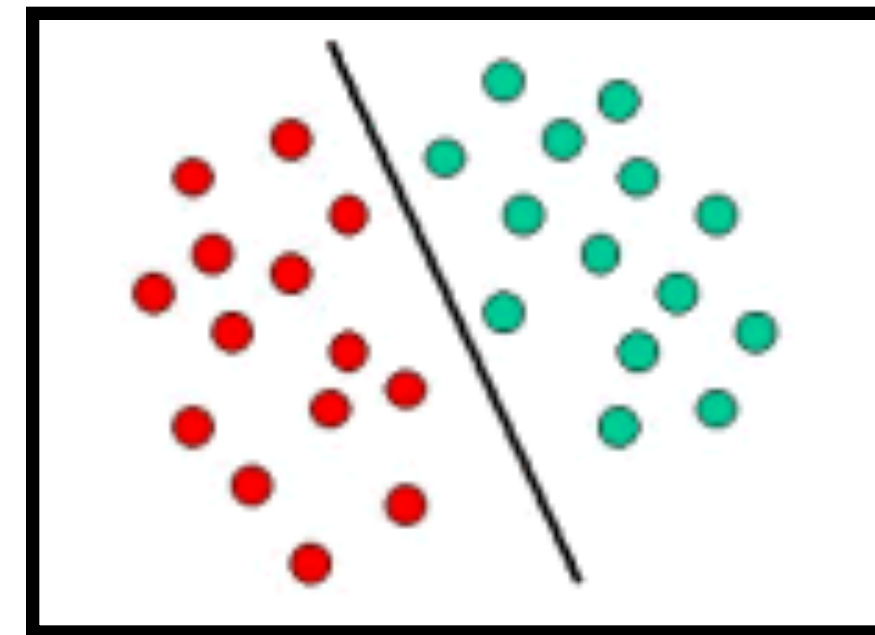
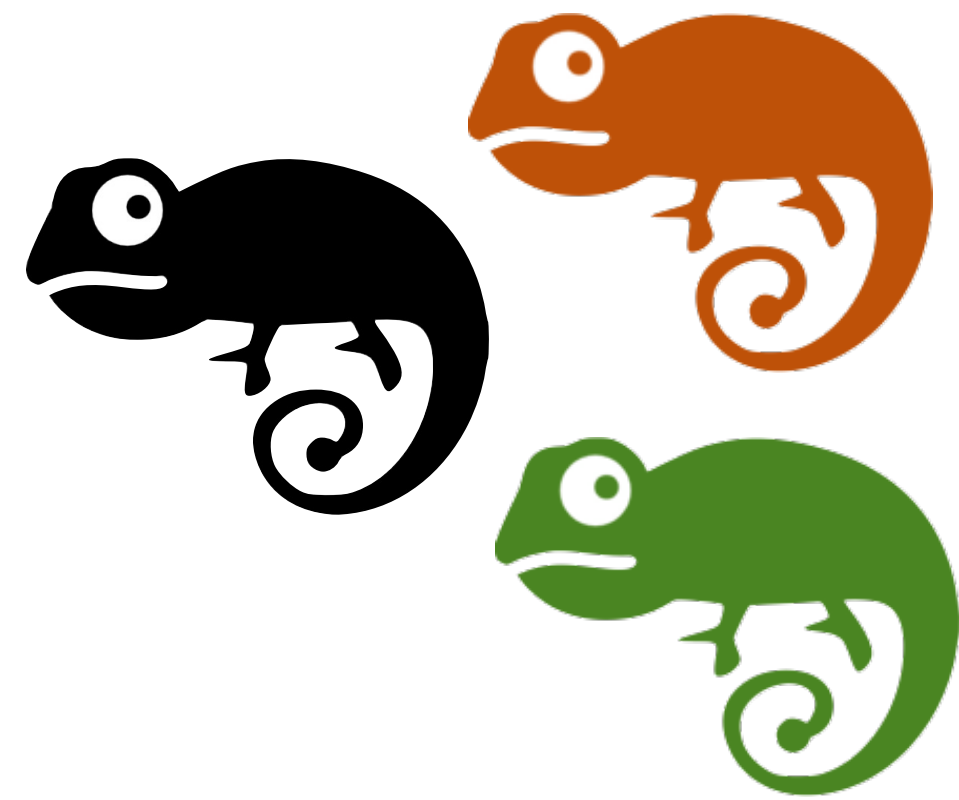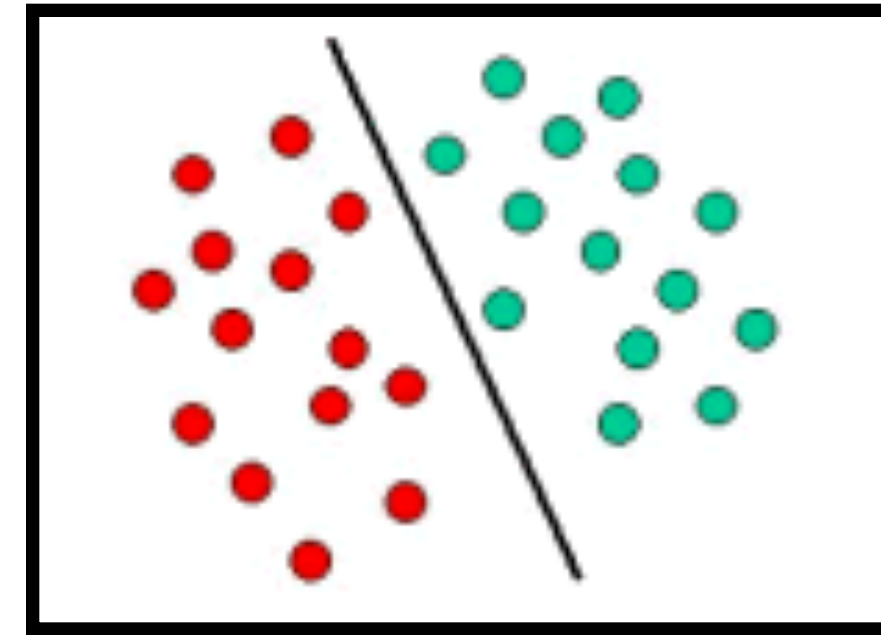Evasion rate on Gmail : 179/380 (47.1%)

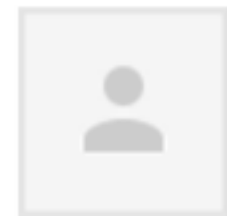# Conclusion



Vs.

Who will win this arm race?

Source Code: http://www.EvadeML.org

# Conclusion



Vs.

# Who will win this arm race?

Source Code: http://www.EvadeML.org

**Ad**: Weilin is seeking summer internship opportunities.

# Backups

# They Don't Care

security@google.com

to me

Hey!

Thanks for your feedback. I think generally because of the ways anti-viruses work there's not really much we can do in this case, but thanks for letting us know!

Eduardo
Google Security Team