# Malware Lineage in the Wild

Irfan Ul Haq*, Sergio Chica*, Somesh Jha[+], and Juan Caballero*

*IMDEA Software Institute, Spain
[+]University of Wisconsin, USA

presented at IMDEA Siminar Series (July 18, 2017)

**Under Submission**

# Lineage

# Lineage

- Software evolution over time
  - New features
  - Bug fixes
  - ...



1985   1992   1995   2001   2006   2009   2015

# Lineage

– Software evolution over time
  – New features
  – Bug fixes
  – ...



1985   1992   1995   2001   2006   2009   2015

– Malware evolves similarly
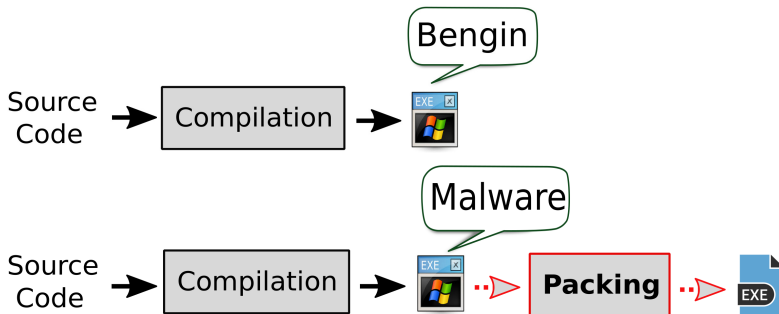
# Background and Motivation

– Unknown versions

# Background and Motivation

- Unknown versions
- Development

# Background and Motivation

- Unknown versions
- Development

# Challenges

– Unknown versions

# Challenges

- Unknown versions
- Unpacking

# Challenges

- Unknown versions
- Unpacking
- Disassembly

# Challenges

- Unknown versions
- Unpacking
- Disassembly
- Granularity

# Challenges

- Unknown versions
- Unpacking
- Disassembly
- Granularity
- Unknown development model

# Challenges

- Unknown versions
- Unpacking
- Disassembly
- Granularity
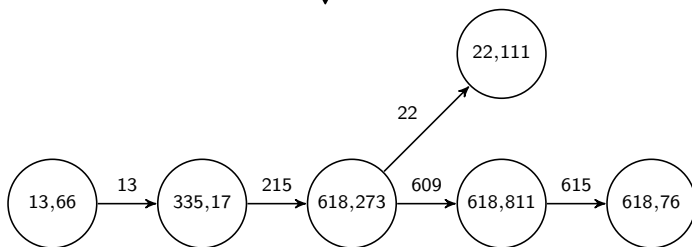- Unknown development model
- Incomplete data
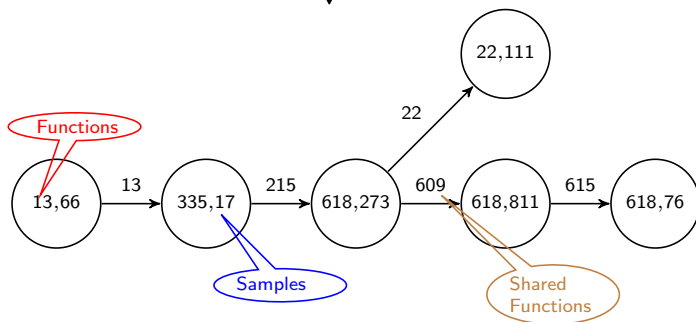
# Why Lineage?

Malware analysis applications

- Triage
- Labeling
- Threat Intelligence
- Author attribution
- ...

# Problem Definition

# Problem Definition

# Problem Definition

## Phase1 → Code Recovery

1. Unpacking
2. Disassembly

# Approach → Overview

## PhaseI → Code Recovery

1. Unpacking
2. Disassembly

## PhaseII → Lineage Inference

1. Identifying Version
2. Building a Lineage Tree
3. Adding Cross-Edges

# Approach → Overview

## PhaseI → Code Recovery

1. Unpacking
2. Disassembly

## PhaseII → Lineage Inference

1. Identifying Version
2. Building a Lineage Tree
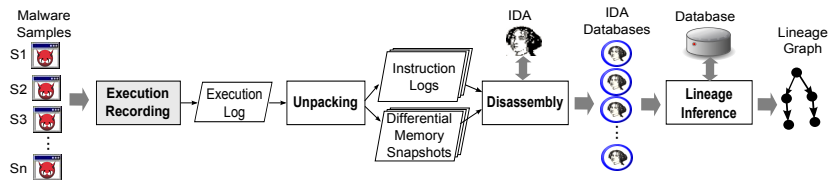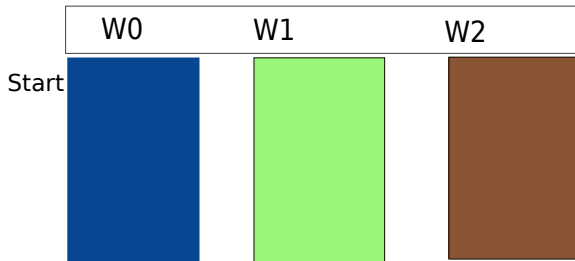3. Adding Cross-Edges

# Code Recovery → Unpacking

# Code Recovery → Unpacking
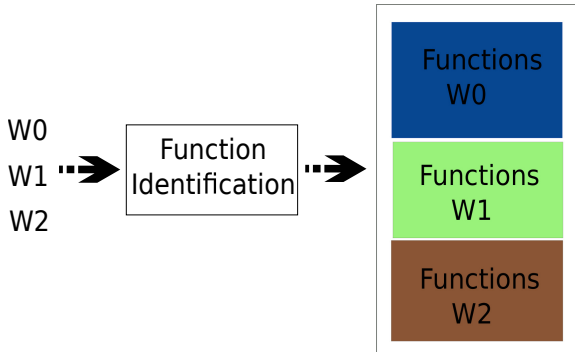
– Monitor each write and execute in memory

# Code Recovery → Unpacking

&ndash; Monitor each write and execute in memory

# Code Recovery → Disassembly

# Lineage Algorithm→ Identifying Versions

– Same functions represent single version

# Lineage Algorithm→ Identifying Versions

- – Same functions represent single version

- – BinDiff (Scalability?)

# Lineage Algorithm→ Identifying Versions

- – Same functions represent single version

- – BinDiff (Scalability?)

- – Function hash (RAW, SPP)

# Lineage Algorithm→ Identifying Versions

- Same functions represent single version

- BinDiff (Scalability?)

- Function hash (RAW, SPP)

$ProgramHash = md5sum(hash(f_1) + hash(f_2) + ... + hash(f_n))$

# Lineage Algorithm→ Buiding Lineage Tree

|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| **P13**    | -   | 12  | 13   | 13   | 13     | 13     |
| **P22**    | -   | -   | 15   | 13   | 22     | 13     |
| **P335**   | -   | -   | -    | 214  | 215    | 214    |
| **P618**   | -   | -   | -    | -    | 609    | 615    |
| **P618_1** | -   | -   | -    | -    | -      | 609    |
| **P618_2** | -   | -   | -    | -    | -      | -      |

# Lineage Algorithm→ Buiding Lineage Tree

|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| P13    | -   | 12  | 13   | 13   | 13     | 13     |
| P22    | -   | -   | 15   | 13   | 22     | 13     |
| P335   | -   | -   | -    | 214  | 215    | 214    |
| P618   | -   | -   | -    | -    | 609    | 615    |
| P618_1 | -   | -   | -    | -    | -      | 609    |
| P618_2 | -   | -   | -    | -    | -      | -      |

Identifying root $\rightarrow$ Lehman's $6th$ law of continuous growth.
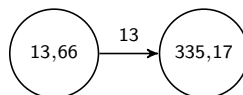
$$\boxed{13,66}$$

# Lineage Algorithm→ Buiding Lineage Tree

|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| **P13**    | -   | 12  | 13   | 13   | 13     | 13     |
| **P22**    | -   | -   | 15   | 13   | 22     | 13     |
| **P335**   | -   | -   | -    | 214  | 215    | 214    |
| **P618**   | -   | -   | -    | -    | 609    | 615    |
| **P618_1** | -   | -   | -    | -    | -      | 609    |
| **P618_2** | -   | -   | -    | -    | -      | -      |

Selected *Nodes* → $P335, P618, P618\_1, P618\_2$
**Break tie:** minimum number of functions

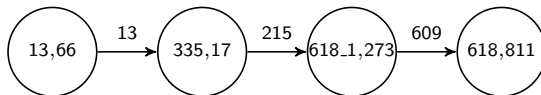|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| P13    | -   | 12  | 13   | 13   | 13     | 13     |
| P22    | -   | -   | 15   | 13   | 22     | 13     |
| P335   | -   | -   | -    | 214  | 215    | 214    |
| P618   | -   | -   | -    | -    | 609    | 615    |
| P618_1 | -   | -   | -    | -    | -      | 609    |
| P618_2 | -   | -   | -    | -    | -      | -      |

Selected *Node* → *P*618_1

# Lineage Algorithm→ Buiding Lineage Tree

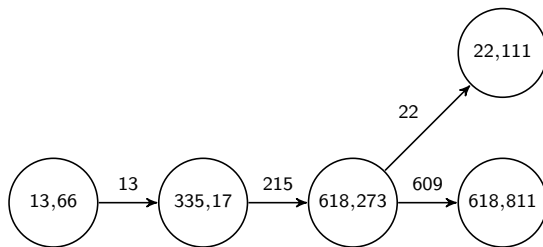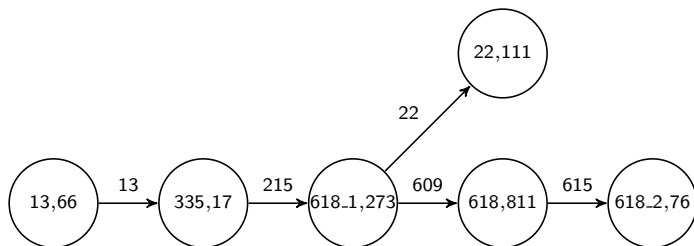|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| P13    | -   | 12  | 13   | 13   | 13     | 13     |
| P22    | -   | -   | 15   | 13   | 22     | 13     |
| P335   | -   | -   | -    | 214  | 215    | 214    |
| P618   | -   | -   | -    | -    | 609    | 615    |
| P618_1 | -   | -   | -    | -    | -      | 609    |
| P618_2 | -   | -   | -    | -    | -      | -      |

Selected *Node* $\rightarrow$ *P*618

# Lineage Algorithm→ Buiding Lineage Tree

|        | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|--------|-----|-----|------|------|--------|--------|
| **P13**    | -   | 12  | 13   | 13   | 13     | 13     |
| **P22**    | -   | -   | 15   | 13   | 22     | 13     |
| **P335**   | -   | -   | -    | 214  | 215    | 214    |
| **P618**   | -   | -   | -    | -    | 609    | 615    |
| **P618_1** | -   | -   | -    | -    | -      | 609    |
| **P618_2** | -   | -   | -    | -    | -      | -      |

Selected *Node* → *P*22

# Lineage Algorithm→ Buiding Lineage Tree

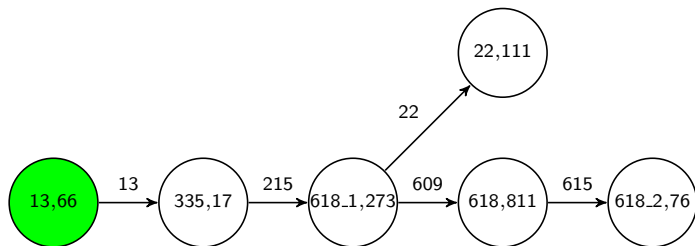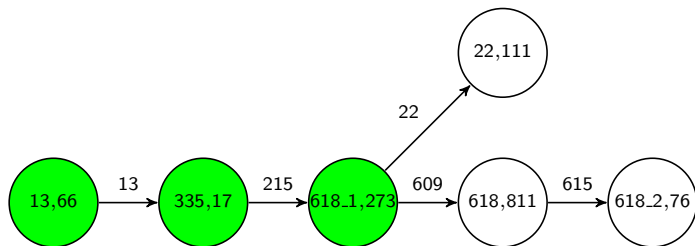|  | P13 | P22 | P335 | P618 | P618_1 | P618_2 |
|---|---|---|---|---|---|---|
| P13 | - | 12 | 13 | 13 | 13 | 13 |
| P22 | - | - | 15 | 13 | 22 | 13 |
| P335 | - | - | - | 214 | 215 | 214 |
| P618 | - | - | - | - | 609 | 615 |
| P618_1 | - | - | - | - | - | 609 |
| P618_2 | - | - | - | - | - | - |

Selected $Node \rightarrow P618\_2$

# Lineage Algorithm→ Adding Cross-Edges

    – Multiple parents, i.e., branching and merging

    – Iterate over each node in topological order

    – Ignore successor and predessors

    – Non-inherited functions

# Lineage Algorithm→ Adding Cross-Edges

- – Multiple parents, i.e., branching and merging
- – Iterate over each node in topological order
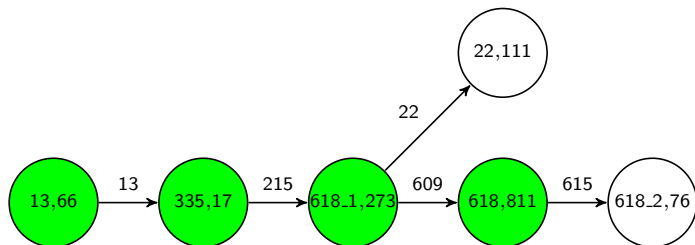- – Ignore successor and predessors
- – Non-inherited functions

# Lineage Algorithm→ Adding Cross-Edges

- Multiple parents, i.e., branching and merging
- Iterate over each node in topological order
- Ignore successor and predessors
- Non-inherited functions

# Lineage Algorithm→ Adding Cross-Edges

- – Multiple parents, i.e., branching and merging
- – Iterate over each node in topological order
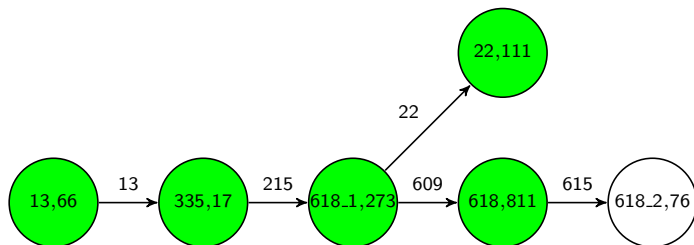- – Ignore successor and predessors
- – Non-inherited functions

# Lineage Algorithm→ Adding Cross-Edges

- Multiple parents, i.e., branching and merging
- Iterate over each node in topological order
- Ignore successor and predessors
- Non-inherited functions

# Evaluation → Goals

- Benign Lineage
- Malware Lineage

# Evaluation → Dataset

## Benign (631 Versions, 13 Programs)

- FileZilla → **10 years**
- Fzputtygen → **10 years**
- Fzsftp → **10 years**
- Notepad++ → **14 years**
- Pageant → **16 years**
- Plink → **16 years**
- ProcessHacker → **6 years**
- PSCS → **16 years**
- PSFTP → **14 years**
- PuTTY → **14 years**
- PuTTYgen → **16 years**
- PuTTYtel → **14 years**
- WinSCP → **10 years**

## Malware (7,793 Samples, 10 Families)

- Allaple → **4,000**
- IRCBot → **365**
- Klez → **750**
- Loring → **216**
- Memery → **113**
- Picsys → **131**
- Simbot → **214**
- Sytro → **1,354**
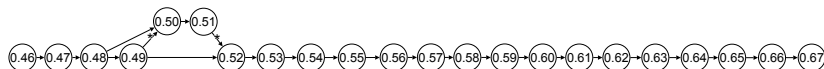- Urelas → **206**
- VtFlooder → **444**

# Evaluation → Benign

Lineage Inference evaluation on open-source programs.

| Program | Reference Type | $|V|$ | SPP $|V|$ | $|V^r|$ | $|X|$ | PO | Raw $|V|$ | $|V^r|$ | $|X|$ | PO |
|---|---|---|---|---|---|---|---|---|---|---|
| FileZilla | S | 119 | 117 | 1 | 8 | 96% | 119 | 1 | 20 | 97% |
| Fzputtygen | S | 107 | 19 | 1 | 0 | 72% | 30 | 1 | 0 | 96% |
| Fzsftp | S | 116 | 50 | 1 | 0 | 71% | 52 | 1 | 2 | 33% |
| Notepad++ | D | 70 | 70 | 1 | 6 | 65% | 70 | 1 | 14 | 71% |
| Pageant | S | 18 | 18 | 1 | 0 | 93% | 18 | 1 | 0 | 97% |
| Plink | S | 18 | 18 | 1 | 0 | 99% | 18 | 1 | 1 | 90% |
| ProcessHacker | 2-S | 52 | 52 | 2 | 2 | 76% | 52 | 2 | 7 | 79% |
| PSCP | S | 20 | 20 | 1 | 1 | 99% | 20 | 1 | 1 | 73% |
| PSFTP | S | 16 | 16 | 1 | 0 | 99% | 16 | 1 | 0 | 88% |
| PuTTY | S | 22 | 22 | 1 | 2 | 100% | 22 | 1 | 0 | 92% |
| PuTTYgen | S | 17 | 16 | 1 | 0 | 87% | 16 | 1 | 0 | 66% |
| PuTTYtel | S | 4 | 4 | 1 | 1 | 100% | 4 | 1 | 0 | 50% |
| WinSCP | S | 47 | 47 | 1 | 1 | 100% | 47 | 1 | 9 | 100% |

# Lineage → Benign

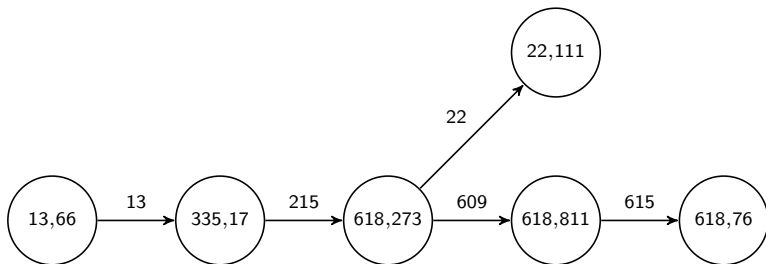**PuTTY**



**Fzputtygen**

# Evaluation → Malware

Lineage Inference evaluation on malware samples.

| Family | EXE | $\|V\|$ | | $\|E\|$ | | $max(\|L(v_i)\|)$ | | $\|L(v_i)=1\|$ | | $max(\|F(v_i)\|)$ | | $min(\|F(v_i)\|)$ | | $\|\bigcup F(v_i)\|$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | spp | raw | spp | raw | spp | raw | spp | raw | spp | raw | spp | raw | spp | raw |
| allaple | 241 | 114 | 143 | 113 | 142 | 71 | 71 | 92 | 138 | 250 | 301 | 10 | 10 | 413 | 510 |
| klez | 118 | 13 | 13 | 12 | 12 | 93 | 93 | 8 | 8 | 618 | 667 | 244 | 251 | 666 | 927 |
| memery | 140 | 12 | 12 | 11 | 11 | 66 | 66 | 3 | 3 | 121 | 123 | 25 | 27 | 131 | 136 |
| picsys | 113 | 5 | 5 | 4 | 4 | 92 | 92 | 1 | 1 | 379 | 473 | 16 | 16 | 397 | 498 |
| simbot | 135 | 21 | 93 | 20 | 92 | 65 | 24 | 13 | 82 | 67 | 71 | 17 | 17 | 108 | 1,723 |
| sytro | 186 | 4 | 4 | 3 | 3 | 92 | 92 | 0 | 0 | 617 | 667 | 335 | 350 | 754 | 1,290 |
| vtflooder | 170 | 20 | 69 | 18 | 66 | 75 | 75 | 11 | 61 | 712 | 749 | 10 | 10 | 3,202 | 4,652 |

# Lineage → Malware



**Picsys**

**Sytro**

# Limitations

– Packers that modify the original code, e.g., VMProtect

– Evasion, anti-VM checks

– Code semantics require manual analysis
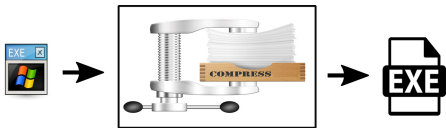
– Function identification, e.g., Nucleus

# Conclusion

– An approach to perform malware lineage on samples collected in the wild.

– First study to identify number of versions in a malware family.

– Evaluated on 13 benign programs and 10 malware families.

– Our approach provides, on average, 26x reduction from input sample versions.

# Thank You!

# Questions?

# Packing/Unpacking

# Packing/Unpacking

# Packing/Unpacking