*Symantec Research Labs*

# Automatic Generation of String Signatures for Malware Detection

*Scott Schneider, Kent Griffin* – SRL
*Xin Hu* – University of Michigan, Ann Arbor
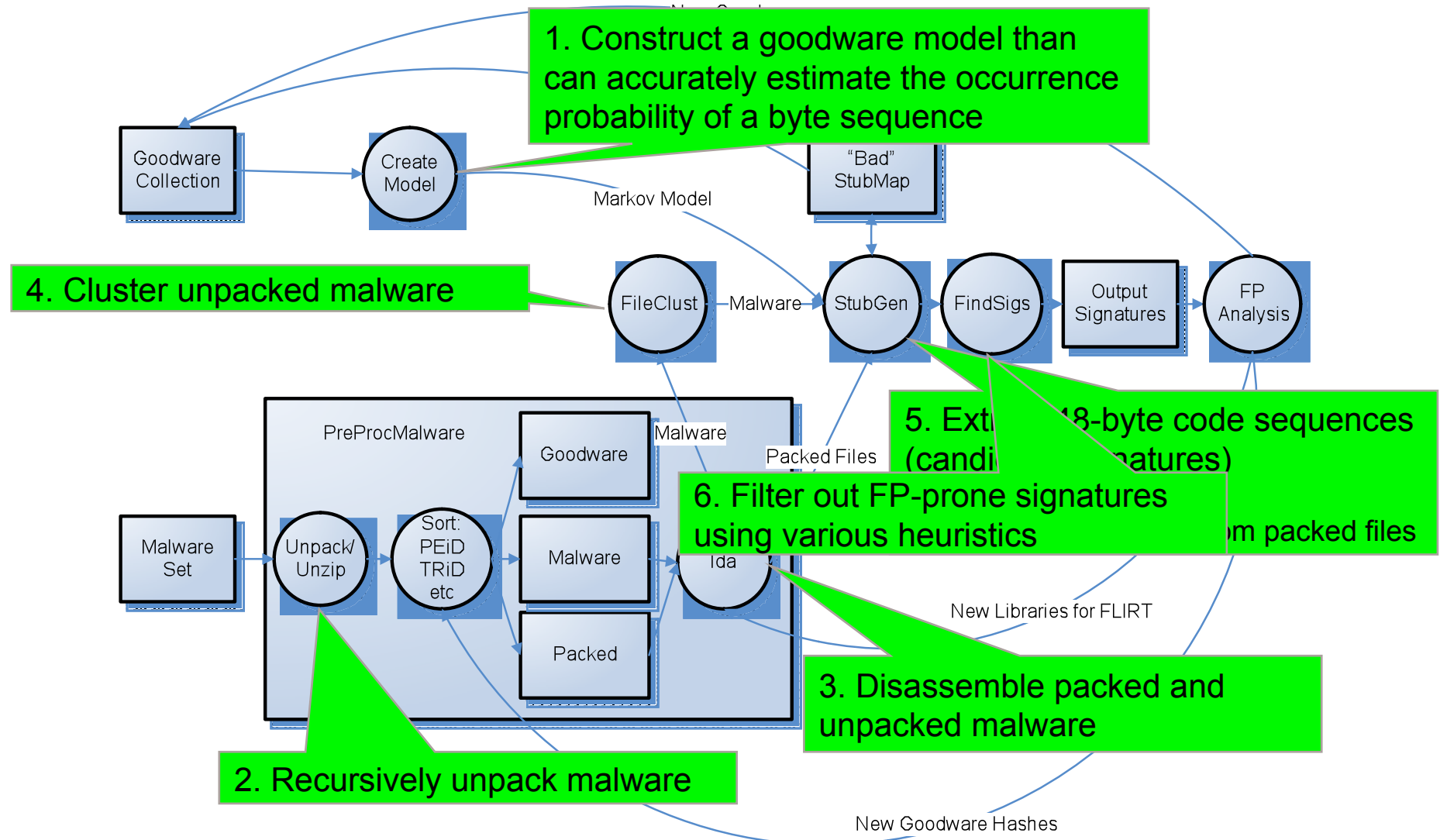*Tzi-cker Chiueh* – Stonybrook University

*September 24, 2009*

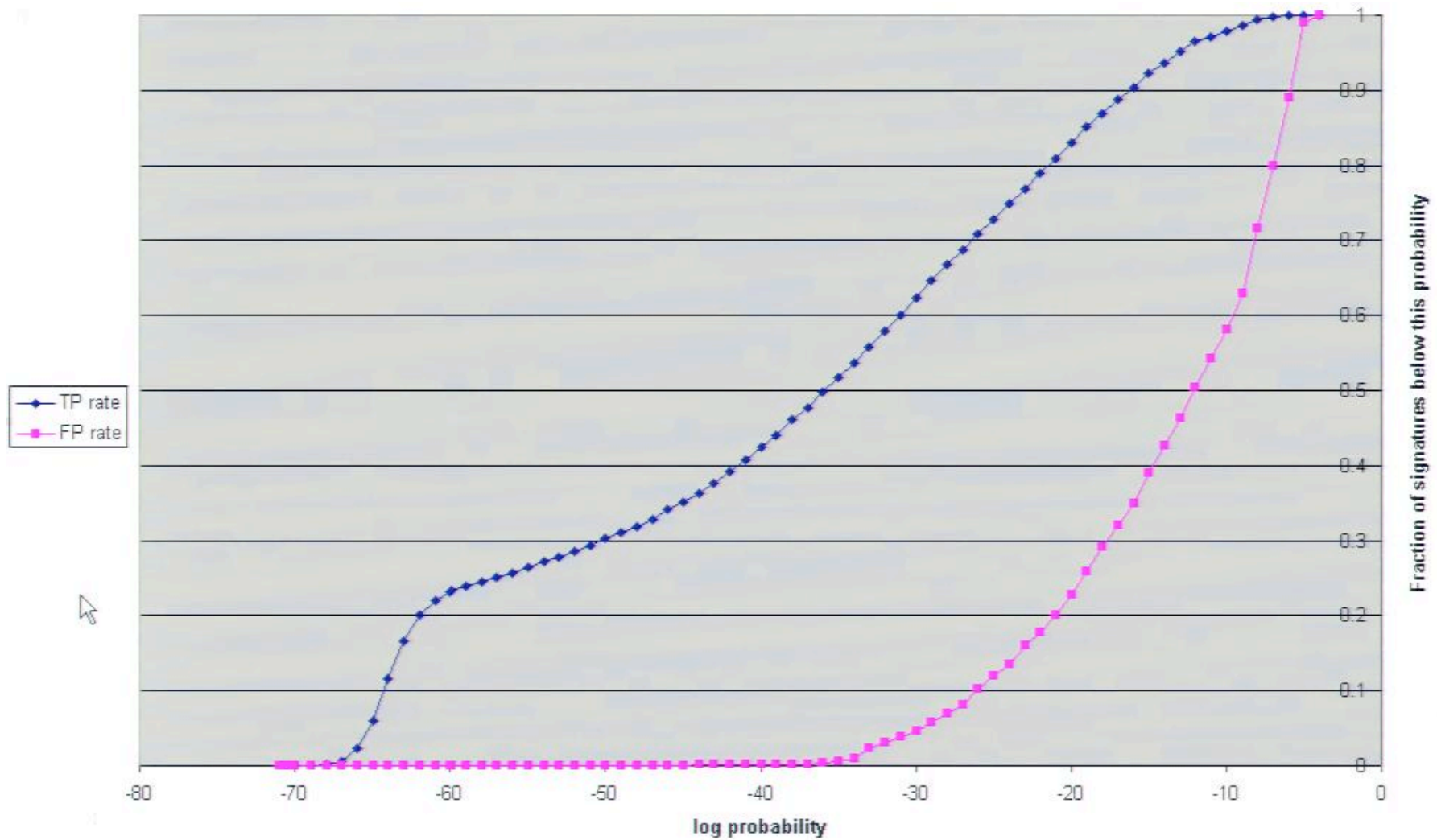# String Signature Generation

- Goal: Given a set of malware samples, derive a minimal set of string signatures that can cover as many malware samples as possible while keeping the FP rate close to zero
  - 48-byte sequences from code

- Why string signatures?
  - Still one of the main techniques for Symantec and other AV companies
  - Higher coverage than file hashes → smaller signature set
  - Currently created manually!

# System Overview

Goodware Collection

Create Model

"Bad" StubMap

**1. Construct a goodware model than can accurately estimate the occurrence probability of a byte sequence**

Markov Model

**4. Cluster unpacked malware**

FileClust — Malware — StubGen — FindSigs — Output Signatures — FP Analysis

**5. Ext... 18-byte code sequences (candi... ...natures)**

PreProcMalware

Goodware

Malware

Packed Files

**6. Filter out FP-prone signatures using various heuristics**

...m packed files

Malware Set

Unpack/ Unzip

Sort: PEiD TRiD etc

Malware

Ida

New Libraries for FLIRT

Packed

**3. Disassemble packed and unpacked malware**

**2. Recursively unpack malware**

New Goodware Hashes

# Heuristics

- 3 main categories:

  - Probability-based – using a Markov chain model

  - Diversity-based – identifies rare libraries and other reused code

  - Disassembly-based – examines assembly instructions

- Discrimination power

  - The best heuristics have high FP reduction and low coverage reduction

  - $\log (FP_i / FP_f) / \log (Coverage_i / Coverage_f)$

  - Raw vs marginal discrimination power

# Goodware Model Effectiveness

# Modeling

- Fixed 5-gram Markov chain model
  - Fixed because the rarest byte sequences are the most important
    - LZ-based training backfired
    - Variable-order models use much more memory

- Needed ~100 MB of relevant data to work

- Probability calculated as in Prediction by Partial Matching
  - p(c|ab) = [c(abc) / c(ab)] * (1-ε(c(ab)))  + p(c|b) * ε(c(ab))
  - ε(c) = sqrt(32) / (sqrt(32) + sqrt(c))

# Scaling the Model

- We have TBytes of training data

  – A model trained on this would use too much memory

  – Solution:  create several models, then prune and merge them

- Pruning

  – If p(c|ab) is close to p(c|b), we don't need node abc

  – If |log(p(c|ab)) – log(p(c|b))| < log(threshold), remove abc

    • Thresholds up to 200 preserve most of the model's effectiveness

# Pruned Model Results



100 MB training data (for pruned case)

# Pruned Model Results Continued



1 GB training data (for pruned case)

# Diversity-based Heuristics

- High coverage signatures are more likely to be from rare library code

  - Model-only tests had 25-30% FPs

- So we examine the *diversity* of covered malware files

  - If files are from many malware families, it's probably a library

# Byte-level Diversity-based Heuristics

- ## Group count/ratio

  - Cluster malware into families

  - Reject signatures that cover too many groups
    or have too high a ratio of groups to covered files

- ## Signature position deviation

  - How much does the signature's position in the files vary?

- ## Multiple common signatures

  - Find a 2$^{nd}$ signature a fixed distance (≥1kb) away in all covered files

# Instruction-level Diversity-based Heuristics

- Enclosing function count
  - Different enclosing functions indicates code reuse
- Several ways of comparing enclosing functions:
  - Exact byte sequences
  - Instruction op codes with some canonicalization
    - e.g. All ADD instructions are treated the same
  - Instruction sequence de-obfuscation
    - e.g. "test esi, esi" and "or esi, esi" is the same

| Method | % FP sig.s Remaining | % all sig.s Remaining | Discrimination Power |
|---|---|---|---|
| Exact byte sequences | 17% | 54% | 2.9 |
| Op code canonicalization | 78% | 90.5% | 2.5 |
| Instruction de-obfuscation | 89% | 94.7% | 2.1 |

# Disassembly-based Heuristics

- IDA Pro's FLIRT –
  Fast Library Identification and Recognition Technology

  – Universal FLIRT

  – Library function reference heuristic

  – Address space heuristic

- Code interestingness…

# Code Interestingness Heuristic

- Encodes Symantec analysts' intuitions using fuzzy logic

- Targets code that is suspicious and/or unlikely to FP

- Points for

  – Unusual constant values

  – Unusual address offsets

    • May indicate custom structs/classes

  – Local, non-library function calls

  – Math instructions

    • Often done by malware for obfuscation

# Results

| Thresholds | Coverage | # sigs | # FPs | # Good sigs | # So-so sigs | # Bad sigs |
|---|---|---|---|---|---|---|
| Loose | 15.7% | 23 | 0 | 6 | 7 | 1 |
| Normal | 14.0% | 18 | 0 | 6 | 2 | 0 |
| Strict | 11.7% | 11 | 0 | 6 | 0 | 0 |
| All non-FP | 22.6% | 220 | 0 | 10 | 11 | 9 |

- Used samples for August 2008
  - 2,363 unpacked files

| Threshold settings | Prob. | Group ratio | Pos. dev. | # common sig.s | Interesting score | Min. coverage |
|---|---|---|---|---|---|---|
| Loose | -90 | 0.35 | 4000 | Single | 13 | 3 |
| Normal | -90 | 0.35 | 3000 | Single | 14 | 4 |
| Strict | -90 | 0.35 | 3000 | Dual | 17 | 4 |

# Results

- 2007-8 files
  - 46,988 unpacked files

| Thresholds | Coverage | # sigs | # FPs |
|---|---|---|---|
| Loose | 14.1% | 1650 | 7 |
| Normal | 11.7% | 767 | 2 |
| Normal + pos. dev. 1,000 | 11.3% | 715 | 0 |
| Strict | 4.4% | 206 | 0 |
| All non-FP | 31.8% | 7305 | 0 |

# Raw Discrimination Power

| Heuristic | % FPs Remaining | % Coverage | Discrimination Power |
|---|---|---|---|
| Position deviation (from ∞ to 8,000) | 41.7% | 96.6% | 25 |
| Min File Coverage (from 3 to 4) | 6.0% | 83.3% | 15 |
| Group Ratio (from 1.0 to .6) | 2.4% | 74.0% | 12 |
| *Probability (from -80 to -100) | 51.2% | 73.7% | 2.2 |
| *Interestingness (from 13 to 15) | 58.3% | 78.2% | 2.2 |
| Multiple common sig.s (from 1 to 2) | 91.7% | 70.2% | 0.2 |
| *Universal FLIRT | 33.1% | 71.7% | 3.3 |
| *Library function reference | 46.4% | 75.7% | 2.8 |
| *Address space | 30.4% | 70.8% | 3.5 |

*Not entirely raw

# Marginal Discrimination Power

| Heuristic | # FPs | % Coverage |
|---|---|---|
| Position deviation (from 3,000 to ∞) | 10 | 121% |
| Min File Coverage (from 4 to 3) | 2 | 126% |
| Group Ratio (from 0.35 to 1) | 16 | 162% |
| Probability (from -90 to -80) | 1 | 123% |
| Interestingness (from 17 to 13) | 2 | 226% |
| Multiple common sig.s (from 2 to 1) | 0 | 189% |
| Universal FLIRT | 3 | 106% |
| Library function reference | 4 | 108% |
| Address space | 3 | 109% |

# Multi-component Signatures

| # Components | # Allowed FPs | Coverage | # Signatures | # FPs |
|:---:|:---:|:---:|:---:|:---:|
| 2 | 1 | 28.9% | 76 | 7 |
| 2 | 0 | 23.3% | 52 | 2 |
| 3 | 1 | 26.9% | 62 | 1 |
| 3 | 0 | 24.2% | 44 | 0 |
| 4 | 1 | 26.2% | 54 | 0 |
| 4 | 0 | 18.1% | 43 | 0 |
| 5 | 1 | 26.2% | 54 | 0 |
| 5 | 0 | 17.9% | 43 | 0 |
| 6 | 1 | 25.9% | 51 | 0 |
| 6 | 0 | 17.6% | 41 | 0 |

- 16 bytes per component, from code and data
- Tested against a smaller goodware set

# Thank You!

Tzi-cker Chiueh

chiueh@cs.sunysb.edu

Kent Griffin

kent_griffin@symantec.com

Xin Hu

huxin@eecs.umich.edu

Scott Schneider

scott_schneider@symantec.com

# Good Signature #0



- Uses 16-bit registers
- Several interesting constants
- Covers 73 files in our malware set
- Very low probability (-140)
- High interestingness score (33)
- Perfect diversity scores

# Good Signature #1

**symantec.**

```
IDA View-A

.text:00010BF2                    add      al, [ebx]
.text:00010BF4
.text:00010BF4 loc_10BF4:                              ; CODE XREF: sub_10BDF+3↑j
.text:00010BF4                                         ; sub_10BDF+9↑j
.text:00010BF4                    call     near ptr loc_10BFD+1
.text:00010BF9                    cmp      [ebx+2Fh], ch
.text:00010BFC                    inc      eax
.text:00010BFD
.text:00010BFD loc_10BFD:                              ; CODE XREF: sub_10BDF:loc_10BF4↑p
.text:00010BFD                    xor      al, [ebx+5E5F04C4h]
.text:00010BFD sub_10BDF         endp ; sp-analysis failed
.text:00010BFD
.text:00010C03                    pop      ebx
.text:00010C04                    pop      ebp
.text:00010C05                    retn
.text:00010C06
.text:00010C06 ; =============== S U B R O U T I N E =====================================
.text:00010C06
.text:00010C06 ; Attributes: bp-based frame
.text:00010C06
.text:00010C06 ; void __stdcall DriverReinitializationRoutine(struct _DRIVER_OBJECT *, PV(
.text:00010C06 DriverReinitializationRoutine proc near ; DATA XREF: DriverReinitializatio
.text:00010C06                                         ; sub_10C95+A↓o
.text:00010C06
.text:00010C06 DriverObject      = dword ptr  8
.text:00010C06
.text:00010C06                    push     ebp
.text:00010C07                    mov      ebp, esp
.text:00010C09                    push     ebx
.text:00010C0A                    push     eax
.text:00010C0B                    push     ebx
.text:00010C0C                    pop      ebx
.text:00010C0D                    pop      eax
.text:00010C0E                    push     0F912h
.text:00010C13                    push     22A6h
.text:00010C18                    push     454Dh
.text:00010C1D                    push     9513h
```

- Several constants
- Covers 65 in our malware set
- Interesting-ness score 19
- Perfect diversity scores

# Good Signature #2



- Several constants
- Covers 63 in our malware set
- Interesting-ness score 21
- Perfect diversity scores

# So-so Signature #4



Suspicious constants – multiples of 10,000

This sig and variants cover 50+ files

Interesting-ness score 13

Good group count, std dev, single sig

Eliminated by better threshold

# So-so Signature #50



- 1 interesting constant
- Covers 4 files in our malware set
- Interestingness score 16
- Good diversity scores
- Eliminated by best thresholds

# Bad Signature #16



- Generic logic
- Only 1 interesting 1-byte constant
- Covers 7 files
- Interestingness score 13
- Bad diversity scores