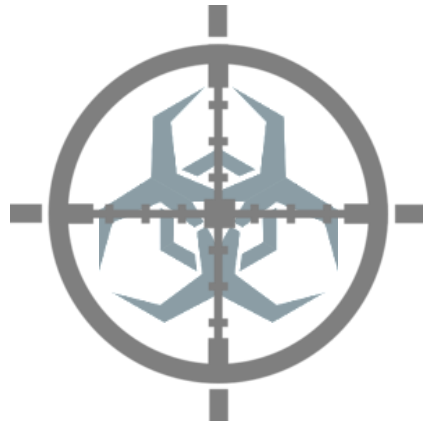# Automated In-memory Malware/Rootkit Detection via Binary Analysis and Machine Learning

By: Malachi Jones, PhD

# ABOUT ME

- **Education**

  - **Bachelors Degree**: Computer Engineering (Univ. of Florida, 2007)

  - **Master's Degree**: Computer Engineering (Georgia Tech, 2009)

  - **PhD**:  Computer Engineering (Georgia Tech, 2013)

- **Cyber Security Experience**

  - **PhD Thesis**: Asymmetric Information Games & Cyber Security  (2009-2013)

  - **Harris Corp.**: Cyber Software Engineer/ Vuln. Researcher (2013-2015)

  - **Booz Allen Dark Labs**: Embedded Security Researcher (2016- Present)

# OUTLINE

- Motivation

- Memory Acquisition Automation

- Phase I Detection: Static Analysis

  i. Efficiently Detecting Dissimilarities in Memory Artifacts

  ii. Clustering Memory Artifacts at Scale to Detect Anomalies

- Phase II Detection: Dynamic Analysis

  i. Automating Dynamic Call Trace Generation

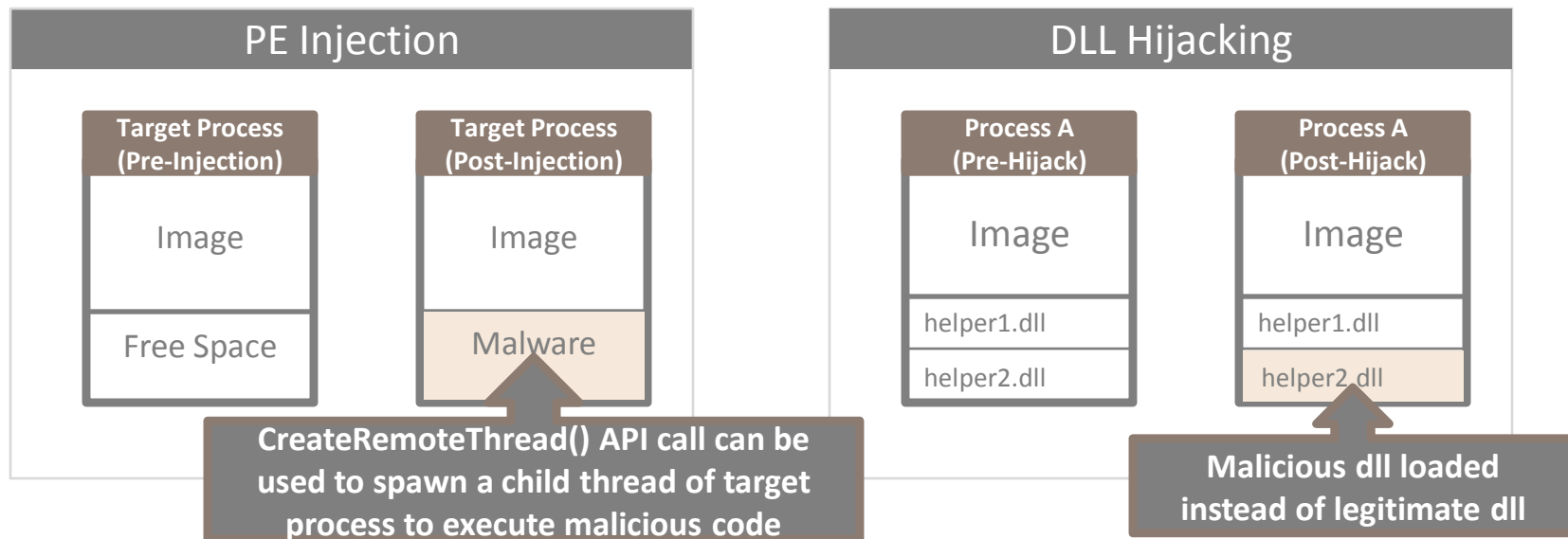  ii. Leveraging Call Trace Information to Detect APTs

- Conclusion/ Q&A

- Appendix

# APPENDIX

- Machine Learning Primer — A

- Advanced Binary Analysis — B

- Memory Forensics — C

- Binary Vector Generation — D

- Call Trace Vector Generation — E

- Hooking — F

# MOTIVATION

- **Code injection** is a technique utilized by malware *to hide malicious code in a legitimate process and/or library* and to force a legitimate process to perform the execution on the malware's behalf

| PE Injection | | DLL Hijacking | |
|---|---|---|---|
| **Target Process (Pre-Injection)** | **Target Process (Post-Injection)** | **Process A (Pre-Hijack)** | **Process A (Post-Hijack)** |
| Image | Image | Image | Image |
| Free Space | Malware | helper1.dll | helper1.dll |
| | | helper2.dll | helper2.dll |

**CreateRemoteThread() API call can be used to spawn a child thread of target process to execute malicious code**

**Malicious dll loaded instead of legitimate dll**

- In addition to **PE Injection** and **DLL Hijacking** *(shown above)*, other methods include **process hollowing** and **reflective DLL Injection**

# MOTIVATION



- **Emulation** and **Hooking** are modern techniques that are employed by anti-virus (AV) vendors *(shown above)* to monitor the execution behavior of binaries executing on a target host

- These techniques combined with **Execution Behavior Analysis** can allow for the discovery of **Advance Persistent Threats** (APT)s that leverage advanced **code injection techniques** to hide in memory and disguise execution

# MOTIVATION

- **Problem**: *Hooking and "traditional" emulation techniques can be reliably evaded by APTs*

- **Examples:**

  - *Hooking -* Malware can either use lower level unhooked APIs or remove hooks at runtime

  - *Emulation-* Utilize an incorrectly implemented API emulation function (e.g. undocumented Windows API) and detect unexpected output given a specified input

# MOTIVATION



- **Observation**: A necessary condition for malicious code to be executed is that the code ***must reside in memory prior to and during execution***

- *As a consequence, periodic live collection and analysis of memory artifacts can provide an effective means to identify malware residing on a host*

# OBJECTIVE

- **Demonstrate** an approach to compare the following *memory artifacts* across a set of networked hosts *in a scalable manner* to identify anomalous code:

    i. Processes

    ii. Shared Libraries

    iii. Kernel Modules

    iv. Drivers

- Specifically, we will discuss how an *approximate clustering algorithm* with *linear run-time performance* can be leveraged to identify outliers among a set of equivalent types of artifacts (e.g. explorer.exe processes) collected from each networked host

- *We will also discuss how **Dynamic Binary Analysis** can be utilized to improve detection of sophisticated malware threats*

# Main Takeaways

## Automated Malware Detection Process

| Static Analysis | → | Cluster | → | Anomaly Detection | → | Dynamic Analysis | → | APT Detection |

**Phase I** (Static Analysis → Cluster → Anomaly Detection)

**Phase II** (Dynamic Analysis → APT Detection)

*Dynamic analysis is leveraged to reduce false positives (from Phase I) and to more accurately identify APTs*

- **Phase I** : We'll leverage **static analysis** to provide us with a computationally efficient way to rapidly identify memory artifacts with anomalous code

- **Phase II**: Dynamic Analysis will be utilized to differentiate between **benign anomalous code** and **malware**

# MAIN TAKEAWAYS

## Set of Networked Hosts

| Host A | Host B | Host C | Host D | Host E | Host F | Host G |

*An agent can be installed on each host to periodically send memory artifacts to a collection server*

## Memory Artifacts Sent to Collection Server

Host A    Host B    Host C    Host D    Host E    Host F    Host G

*An example type of artifact is an explorer.exe process*

## Identical types of artifacts are clustered

Host D    Host B    Host G    |    Host F    |    Host E    Host A    Host C

*Outliers (e.g. Host F) are analyzed dynamically to more accurately identify malicious behavior*

# MAIN TAKEAWAYS

## Set of Networked Hosts

| Host A | Host B | Host C | Host D | Host E | Host F | Host G |

*An agent can be installed on each host to periodically send memory artifacts to a collection server*

## Memory Artifacts Sent to Collection Server

| Host A | Host B | Host C | Host D | Host E | Host F | Host G |

*An example type of artifact is an explorer.exe process*

## Identical types of artifacts are clustered

| Host D | Host B | Host G | | Host F | | Host E | Host A | Host C |

*Outliers (e.g. Host F) are analyzed dynamically to more accurately identify malicious behavior*

# MAIN TAKEAWAYS

## Set of Networked Hosts

Host A   Host B   Host C   Host D   Host E   Host F   Host G

*An agent can be installed on each host to periodically send memory artifacts to a collection server*

## Memory Artifacts Sent to Collection Server

Host A   Host B   Host C   Host D   Host E   Host F   Host G

*An example type of artifact is an explorer.exe process*

## Identical types of artifacts are clustered

Host D   Host B   Host G        Host F        Host E   Host A   Host C

*Outliers (e.g. Host F) are analyzed dynamically to more accurately identify malicious behavior*

# MAIN TAKEAWAYS

## Set of Networked Hosts

Host A  Host B  Host C  Host D  Host E  Host F  Host G

*An agent can be installed on each host to periodically send memory artifacts to a collection server*

## Memory Artifacts Sent to Collection Server

Host A  Host B  Host C  Host D  Host E  Host F  Host G

*An example type of artifact is an explorer.exe process*

## Identical types of artifacts are clustered

Host D  Host B  Host G          Host F          Host E  Host A  Host C

*Outliers (e.g. Host F) are analyzed dynamically to more accurately identify malicious behavior*

# MEMORY ACQUISITION AUTOMATION

# MEMORY ACQUISITION AUTOMATION



| Set of Networked Hosts |
| --- |

| Host A | Host B | Host C |

| Memory Artifacts Sent to Collection Server |
| --- |

*An example of process artifacts acquired from physical memory of each host*

- **Rekall** is an open source tool that can be used to acquire *live* memory and perform analysis

- We can develop an agent that is deployed on each host that interfaces with **Rekall** to collect desired memory artifacts in an automated fashion

# MEMORY ACQUISITION AUTOMATION

- Querying for active and terminated processes



**pslist** *command allows for live querying of internal data structures in memory for active and terminated processes*

- Querying for active and terminated processes



*Terminated process (notepad++) still in memory 26 days later*

`pslist` *command allows for live querying of internal data structures in memory for active and terminated processes*

# MEMORY ACQUISITION AUTOMATION

- Capturing live dumps of binaries w/ Rekall

```
[1] Live (Memory) 16:35:15> procdump pids=[332,324,368,380], dump_dir="C:/tmp/dumps"
----------------------------> procdump(pids=[332,324,368,380], dump_dir="C:/tmp/dumps")
2017-08-30 16:35:15,786:DEBUG:rekall.1:Running plugin (procdump) with args (()) kwargs
umps'})
                      _EPROCESS                          Filename
---------------------------------------------------- --------
0xfa800d9deb10 smss.exe                     324 executable.smss.exe_324.exe
0xfa80106af8e0 notepad++.exe                332 executable.notepad.exe_332.exe
0xfa800ef96060 stacsv64.exe                 368 executable.stacsv64.exe_368.exe
0xfa800d663460 AmazonDrive.ex               380 executable.AmazonDrive.ex_380.exe
Out<16:35:17> Plugin: procdump (ProcExeDump)
[1] Live (Memory) 16:35:17>
```

**procdump** *dumps a set of specified processes (given pids as input) to a desired directory*

# MEMORY ACQUISITION AUTOMATION

- Capturing live dumps of binaries w/ Rekall



**procdump** *dumps a set of specified processes (given pids as input) to a desired directory*

# (DEMO) MEMORY ACQUISITION AUTOMATION

# (DEMO) MEMORY ACQUISITION AUTOMATION

# (Demo) Memory Acquisition Automation

Dumped processes received from client

# (DEMO) MEMORY ACQUISITION AUTOMATION



Process dumps disassembled w/ Binary Ninja

# PHASE I DETECTION: STATIC ANALYSIS

# PHASE I DETECTION: STATIC ANALYSIS

### Memory Artifacts Sent to Collection Server



Host A     Host B     Host C     Host D     Host E     Host F     Host G

### Identical types of artifacts are clustered



Host D     Host B     Host G          Host F          Host E     Host A     Host C

- **Overall Goals**:

  - *Group memory artifacts based on their similarity (as demonstrated in the above figure) to identify outliers*

  - *Leverage Dynamic Analysis (Phase II) to differentiate between benign anomalous code and malware*

# PHASE I DETECTION: STATIC ANALYSIS

| Identical types of artifacts are clustered | | |
|---|---|---|
| Host D   Host B   Host G | Host F | Host E   Host A   Host C |

- **Requirements for clustering artifacts *at scale***

  - **Computationally efficient** method for determining the similarity (or dissimilarity) of a pair of binaries

  - Clustering algorithm with a **linear run-time performance in the worst-case**

Computationally Efficient Diffing Algorithm

# PHASE I DETECTION: STATIC ANALYSIS

# EFFICIENT DIFFING ALGORITHM

- **Question**: Why is efficient binary diffing *critical* to our goal of detecting malicious code?

  - **Slow diffing → Slow clustering → Delayed threat detection**

  - We want to be able to **cluster a large set of binaries (10,000+)** pretty quickly to identify binaries that pose potential threats to hosts on the network

  - Clustering algorithms need to **perform a large number of diffing** operations with respect to the binaries (exact number depends on run-time algorithm performance)

# EFFICIENT DIFFING ALGORITHM

## Step 1: Generate a vector representation of each binary

**main (explorer.exe)**

```
push ebp
mov rbp, rsp
sub rsp, 0x10
```

*Each row of the vector represents the number of occurrences of a unique sequence of instructions*

## Step 2: Compute the similarity of a pair of vectors

explorerA.exe

explorerB.exe

explorerA.exe

explorerB.exe

**Similarity Function**

.867

*Similarity function takes as input two vectors and produces a value between 0 and 1*

**(See Appendix D for more details about the algorithm)**

# EFFICIENT DIFFING ALGORITHM

- Evaluating similarity function of on-disk copy of explorer.exe vs. in-memory image

# EFFICIENT DIFFING ALGORITHM

- Evaluating similarity function against explorer.exe

```
/home/targaryen/.virtualenvs/angr/bin/python /targaryen/targaryen/BinaryReportAnalyzer/ReportSimilarityAnalysis.
INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:
Report Similarity Analysis

INFO:BinaryReportAnalyzer.ReportAnalysis:
(Elasped Time: 0.0248849391937) Deserialized report analysis
    file 'executable.1844.exe'
    hash:'adaca26eb1685da66c547e01363664cc3bf38c2a08a6287044d17690a75bf628'

INFO:BinaryReportAnalyzer.ReportAnalysis:
(Elasped Time: 0.0328030586243) Deserialized report analysis
    file 'explorer_memory_forensics.exe'
    hash:'6bed1a3a956a859ef4420feb2466c040800eaf01ef53214ef9dab53aeff1cff0'

INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Beginning Jacard index analysis
INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Finished Jacard index analysis (Elapsed time:0.0579369068146)
INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Jacard index: 0.834103965252

Process finished with exit code 0
```

**Jaccard Index:
83.4% similar**

# EFFICIENT DIFFING ALGORITHM

- Comparing similarity results against BinDiff

# EFFICIENT DIFFING ALGORITHM

- Comparing similarity results against BinDiff

# Efficient Diffing Algorithm

■ Performance vs BinDiff



```
INFO:BinaryReportAnalyzer.ReportAnalysis:
(Elasped Time: 0.0248849391937) Deserialized report analysis
    file 'executable.1844.exe'
    hash:'adaca26eb1685da66c547e01363664cc3bf38c2a08a6287044d17690a75bf628'

INFO:BinaryReportAnalyzer.ReportAnalysis:
(Elasped Time: 0.0328030586243) Deserialized report analysis
    file 'explorer_memory_forensics.exe'
    hash:'6bed1a3a956a859ef4420feb2466c040800eaf01ef53214ef9dab53aeff1cff0'

INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Beginning Jacard index analysis
INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Finished Jacard index analysis
INFO:BinaryReportAnalyzer.ReportSimilarityAnalysis:Jacard index: 0.834103965252
```

**Jaccard Performance: 57 ms**

(Elapsed time:0.0579369068146)

- **BinDiff Analysis** : 4.2 seconds

- **Jaccard  Index Analysis**: 57 ms  (in pure python)

- *Speed up factor of 74*

# PHASE I DETECTION: STATIC ANALYSIS

- Clustering Algorithms

  - We'll utilize an **agglomerative hierarchical clustering algorithm** because we don't have to specify the exact number of clusters, *k*, a priori (vs. k-means, where *k* must be specified)

  - Computational complexity for a non-approximated implementation of the algorithm is $O(n^2 \log n)$, which is not a desirable property for achieving scalability

  - Instead, we'll use an **approximate implementation** (presented in [1]), which has **a linear worse-case complexity** of $O(k^2 \log k + n)$

  - **Note**: *k is constant and k << n*

# Clustering Artifacts at Scale

- **Approximate Clustering Algorithm [1] Sketch**

| Set of memory artifacts |
|---|
| Host A · Host B · Host C · Host D · Host E · Host F · Host G |

| Step 1: Prototype Extraction $[O(k{\cdot}n)]$ |
|---|
| Host A · Host B · Host F |

*Prototypes are a small (k << n), yet representative subset of artifacts*

| Step 2: Clustering w/ Prototype $[O(k^2 \log(k) + n)]$ |
|---|

**Host A** — Host E, Host C  
**Host B** — Host D, Host G  
**Host F**

■ Prototype Extraction Algorithm [1]

**Algorithm 1** Prototype extraction

1: $prototypes \leftarrow \varnothing$
2: $distance[x] \leftarrow \infty$ for all $x \in reports$
3: **while** $\max(distance) > d_p$ **do**
4:       choose $z$ such that $distance[z] = \max(distance)$
5:       **for** $x \in reports$ and $x \neq z$ **do**
6:           **if** $distance[x] > ||\hat{\varphi}(x) - \hat{\varphi}(z)||$ **then**
7:              $distance[x] \leftarrow ||\hat{\varphi}(x) - \hat{\varphi}(z)||$
8:       add $z$ to $prototypes$

# PHASE II DETECTION: DYNAMIC ANALYSIS

**Automated Malware Detection Process**



- Although dynamic analysis can be more accurate (i.e. fewer false-positives) than static analysis, it can also be more *computationally expensive*

- Therefore, we've discussed utilizing static analysis (Phase I) to filter out binaries based on their similarity.

- Consequently, *we can focus computationally efforts in Phase II on differentiating benign anomalous code from APTs via dynamic analysis*

# PHASE II: DYNAMIC ANALYSIS

- **Key Concept:** *Binaries utilize system calls and library function calls (e.g. dlls) to interact with the operating system in order to perform meaningful/desired operations*

- **Corollary**: *By analyzing the external call sequences of a binary (e.g. call trace), the underlying behavior and intent of the binary can be characterized*

## Case Study: "Sample J" Malware



```
push    ebp
mov     ebp, esp
sub     esp, 130h
push    edi
sidt    fword ptr [ebp+var_8]
mov     eax, dword ptr [ebp+var_8+2]
cmp     eax, 8003F400h
jbe     short loc_10001C88
```

```
cmp     eax, 80047400h
jnb     short loc_10001C88
```

```
xor     eax, eax
pop     edi
mov     esp, ebp
pop     ebp
retn    0Ch
```

```
loc_10001C88:
xor     eax, eax
mov     ecx, 49h
lea     edi, [ebp+pe.cntUsage]
mov     [ebp+pe.dwSize], 0
push    eax            ; th32ProcessID
push    2              ; dwFlags
rep stosd
call    CreateToolhelp32Snapshot
mov     edi, eax
cmp     edi, 0FFFFFFFFh
jnz     short loc_10001CB9
```

```
xor     eax, eax
pop     edi
mov     esp, ebp
pop     ebp
retn    0Ch
```

```
loc_10001CB9:
lea     eax, [ebp+pe]
push    esi
push    eax            ; lppe
push    edi            ; hSnapshot
mov     [ebp+pe.dwSize], 128h
call    Process32First
test    eax, eax
jz      short loc_10001D24
```

```
mov     esi, ds:_stricmp
lea     ecx, [ebp+pe.szExeFile]
```

```c
 1 BOOL __stdcall DllMain(HINSTANCE hinstDLL, DWORD fdwReason,
 2 {
 3   HANDLE v4; // edi@4
 4   DWORD v5; // eax@10
 5   DWORD v6; // ecx@10
 6   PROCESSENTRY32 pe; // [esp+4h] [ebp-130h]@4
 7   char v8[6]; // [esp+12Ch] [ebp-8h]@1
 8
 9   __sidt(v8);
10   if ( *(_DWORD *)&v8[2] > 0x8003F400 && *(_DWORD *)&v8[2] <
11     return 0;
12   pe.dwSize = 0;
13   memset(&pe.cntUsage, 0, 0x124u);
14   v4 = CreateToolhelp32Snapshot(2u, 0);
15   if ( v4 == (HANDLE)-1 )
16     return 0;
17   pe.dwSize = 296;
18   if ( Process32First(v4, &pe) )
19   {
20     if ( !stricmp(pe.szExeFile, Str2) )
21     {
22 LABEL_10:
23       v5 = pe.th32ParentProcessID;
24       v6 = pe.th32ProcessID;
25       goto LABEL_12;
26     }
27     while ( Process32Next(v4, &pe) )
28     {
29       if ( !stricmp(pe.szExeFile, Str2) )
30         goto LABEL_10;
31     }
32   }
33   v5 = fdwReason;
34   v6 = fdwReason;
35 LABEL_12:
36   if ( v5 == v6 )
37     return 0;
38   if ( fdwReason == 1 )
39     CreateThread(0, 0, StartAddress, 0, 0, 0);
40   return 1;
```

(**sha1**: 70cb0b4b8e60dfed949a319a9375fac44168ccbb)

# PHASE II: DYNAMIC ANALYSIS

## Case Study: "Sample J" Malware

```
BOOL __stdcall DllMain(HINSTANCE hinstDLL, DWORD fdwReason,
{
  HANDLE v4; // edi@4
  DWORD v5; // eax@10
  DWORD v6; // ecx@10
  PROCESSENTRY32 pe; // [esp+4h] [ebp-130h]@4
  char v8[6]; // [esp+12Ch] [ebp-8h]@1

  __sidt(v8);
  if ( *(_DWORD *)&v8[2] > 0x8003F400 && *(_DWORD *)&v8[2] <
    return 0;
  pe.dwSize = 0;
  memset(&pe.cntUsage, 0, 0x124u);
  v4 = CreateToolhelp32Snapshot(2u, 0);
  if ( v4 == (HANDLE)-1 )
    return 0;
  pe.dwSize = 296;
  if ( Process32First(v4, &pe) )
  {
    if ( !stricmp(pe.szExeFile, Str2) )
    {
LABEL_10:
      v5 = pe.th32ParentProcessID;
      v6 = pe.th32ProcessID;
      goto LABEL_12;
    }
    while ( Process32Next(v4, &pe) )
    {
      if ( !stricmp(pe.szExeFile, Str2) )
        goto LABEL_10;
    }
  }
  v5 = fdwReason;
  v6 = fdwReason;
LABEL_12:
  if ( v5 == v6 )
    return 0;
  if ( fdwReason == 1 )
    CreateThread(0, 0, StartAddress, 0, 0, 0);
  return 1;
```

### Sample J Call Trace Example

```
memset(,0, 0x124u)

CreateToolhelp32Snapshot()

Process32First()

stricmp(., "explorer.exe")

Process32Next(.,.)

stricmp(., "explorer.exe")

CreateThread(,,StartAddress)
```

# PHASE II: DYNAMIC ANALYSIS

**Case Study: "Sample J" Malware**

| Sample J Call Trace Example |
|---|
| `memset(,0, 0x124u)` |
| `CreateToolhelp32Snapshot()` |
| `Process32First()` |
| `stricmp(., "explorer.exe")` |
| `Process32Next(.,.)` |
| `stricmp(., "explorer.exe")` |
| `CreateThread(,,StartAddress)` |

- **Call Trace Analysis**
  i. Iterate through the process handles to see if a process with name "explorer.exe" exists *(Check if user logged in)*
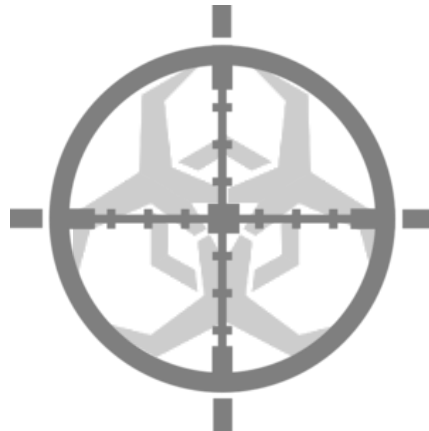  ii. If process exists, create a thread that infects target system

- **Questions:**

    1. How do we automate the call trace generation process in a manner that maximizes traversal of unique code paths?

    2. How do we leverage generated call trace information to identify potential APTs?

# PHASE II  DETECTION

# Automating Call Trace Generation

| 📁 → | Disassembler | → | vex ir | → | ir emulation | → | Call Traces |
|------|--------------|---|--------|---|--------------|---|-------------|

- Call Trace Generation
  i.   **Load** the dumped binary into a disassembler (e.g. IDA or Binary Ninja) and extract the executable portion of binary

  ii.  **Lift** the executable portion of binary to vex, a RISC-like intermediate representation (ir) language

  iii. **Perform** emulation on the vex ir to traverse unique code paths that originate from a specified entry function

  iv.  **Record** calls that occur during traversal of a code path

# AUTOMATING CALL TRACE GENERATION

# (DEMO) CALL TRACE GENERATION

## Target Binary: explorer.exe

# (DEMO) CALL TRACE GENERATION

# (DEMO) CALL TRACE GENERATION

# (Demo) Call Trace Generation

Leverage Call Trace Info to Detect APTs

# PHASE II  DETECTION
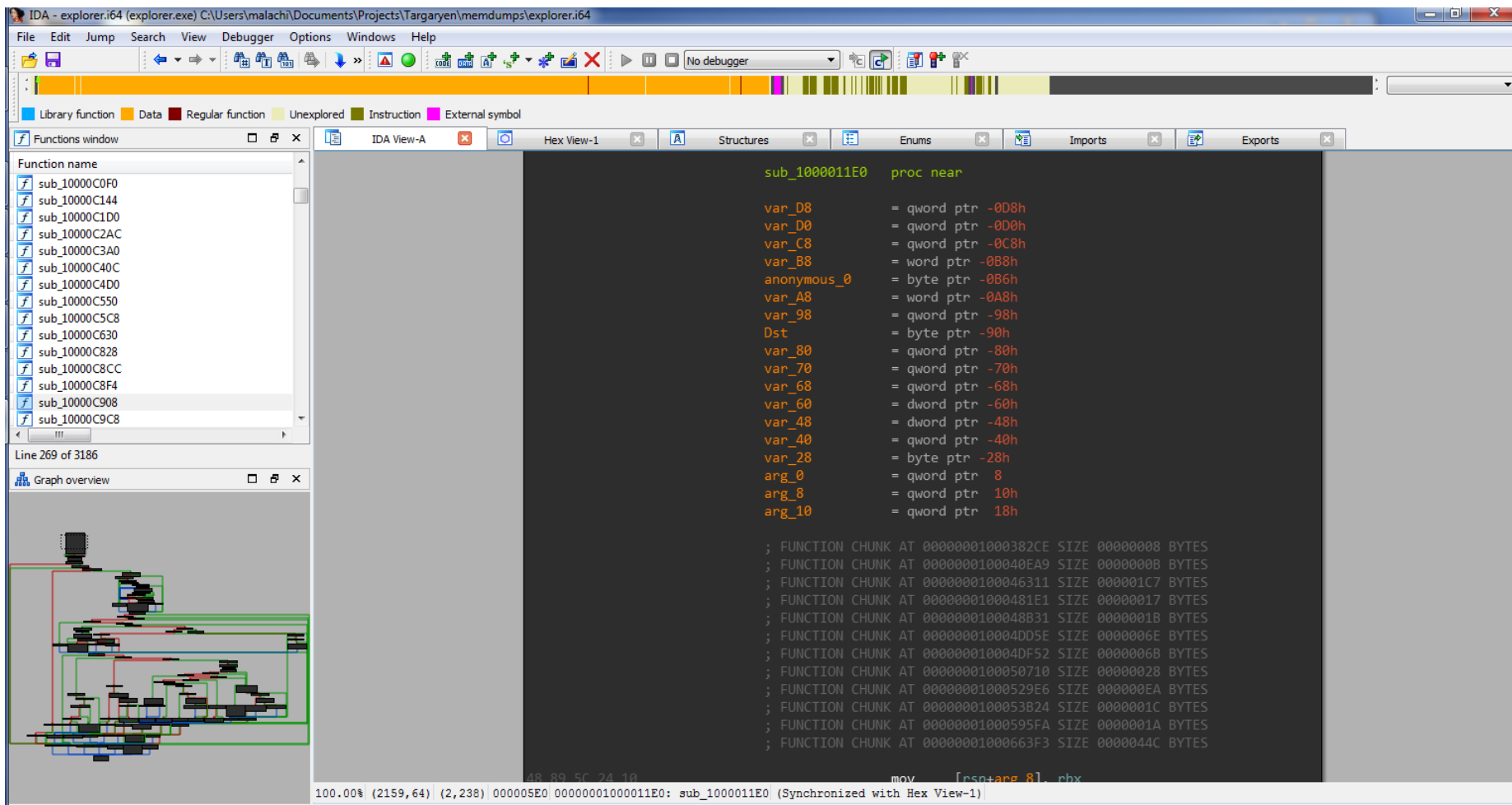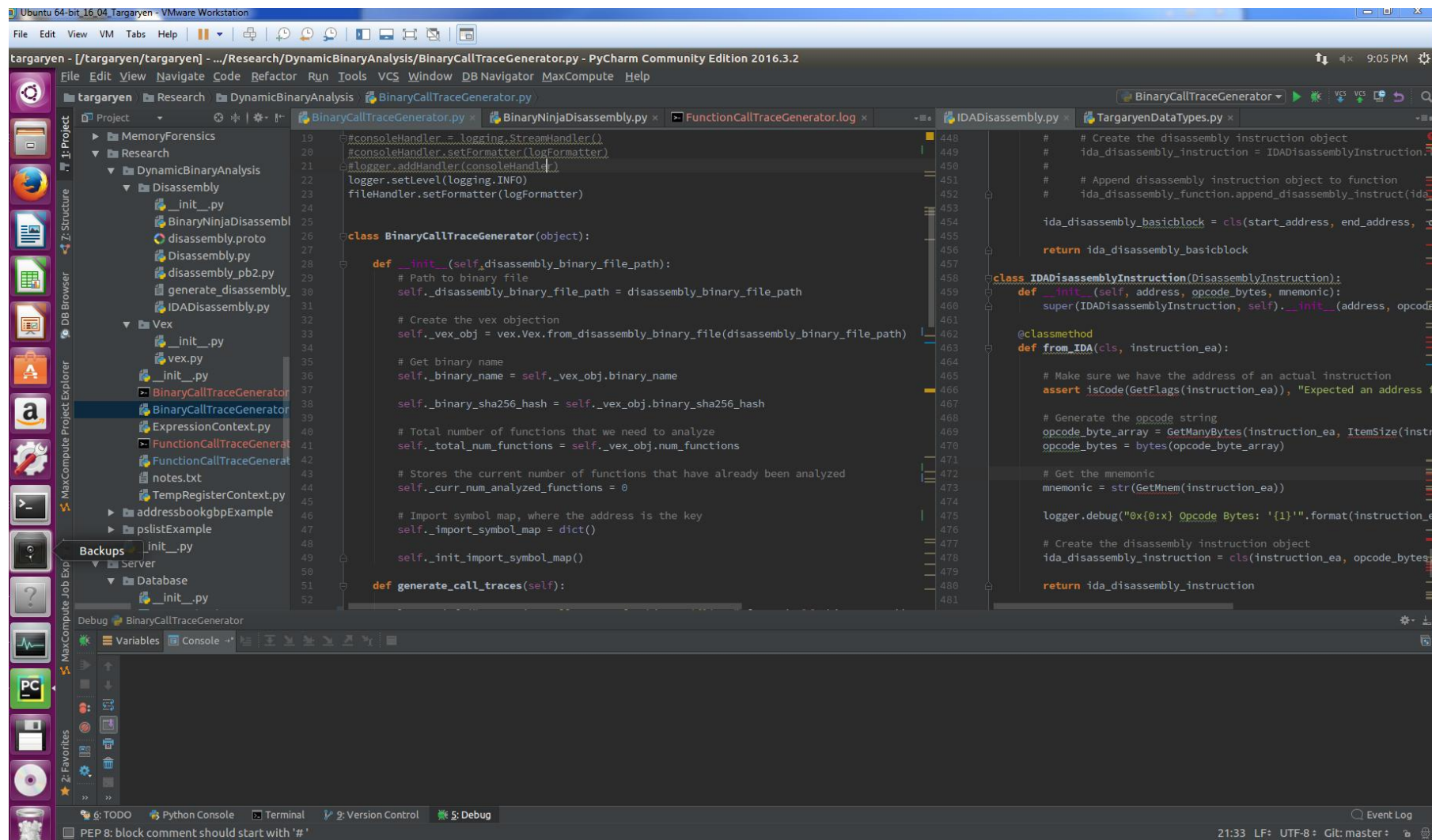
# LEVERAGE CALL TRACE INFO TO DETECT APTs

## Step 1: Generate a vector representation of the set of call traces

**Call Trace 1**
…
…

**Call Trace 2**
…
…

**Call Trace n**
…
…

*Traces generated via **dynamic analysis**, where different code paths are traversed to maximize code coverage*

## Step 2: Build a unified trusted call trace vector from trusted binaries

**Trusted Call Trace vectors**

[ ] + [ ] + [ ] + [ ] + [ ] + [ ]

**Unified call trace vector**

## Step 3: Compare target binary against unified trusted call trace vector

Target vector    Trusted Vector

Target vector

Trusted vector

**Comparison Function**

.6

*Function output is percentage of call sequences in target that are trusted (1 ➔ all sequences trusted)*

# LEVERAGE CALL TRACE INFO TO DETECT APTs

## Step 1: Generate a vector representation of the set of call traces

| Call Trace 1 | Call Trace 2 | Call Trace n |
|---|---|---|
| ... | ... | ... |
| ... | ... | ... |

*Traces generated via **dynamic analysis**, where different code paths are traversed to maximize code coverage*

## Step 2: Build a unified trusted call trace vector from trusted binaries

**Trusted Call Trace vectors**

[ ] + [ ] + [ ] + [ ] + [ ] + [ ]

**Unified call trace vector**

## Step 3: Compare target binary against unified trusted call trace vector

Target vector   Trusted Vector

Target vector

Trusted vector

**Comparison Function**

.6

*Function output is percentage of call sequences in target that are trusted (1 ➔ all sequences trusted)*

# LEVERAGE CALL TRACE INFO TO DETECT APTS

## Step 1: Generate a vector representation of the set of call traces

| Call Trace 1 | Call Trace 2 | Call Trace n |
| --- | --- | --- |
| ... | ... | ... |
| ... | ... | ... |

*Traces generated via **dynamic analysis**, where different code paths are traversed to maximize code coverage*

## Step 2: Build a unified trusted call trace vector from trusted binaries

**Trusted Call Trace vectors**

[ ] + [ ] + [ ] + [ ] + [ ] + [ ]

**Unified call trace vector**

## Step 3: Compare target binary against unified trusted call trace vector

Target vector   Trusted Vector

Target vector

Trusted vector

**Comparison Function**

.6

*Function output is percentage of call sequences in target that are trusted (1 ➜ all sequences trusted)*

# LEVERAGE CALL TRACE INFO TO DETECT APTS

## Step 1: Generate a vector representation of the set of call traces



**Call Trace 1**
...
...

**Call Trace 2**
...
...

**Call Trace n**
...
...

*Traces generated via **dynamic analysis**, where different code paths are traversed to maximize code coverage*

## Step 2: Build a unified trusted call trace vector from trusted binaries

**Trusted Call Trace vectors**

[ ] + [ ] + [ ] + [ ] + [ ] + [ ]

**Unified call trace vector**

## Step 3: Compare target binary against unified trusted call trace vector

Target vector    Trusted Vector

Target vector →

Trusted vector →

**Comparison Function**

.6

*Function output is percentage of call sequences in target that are trusted (1 ➔ all sequences trusted)*

**(See Appendix E for more details about the algorithm)**

# CONCLUSION

# CONCLUSION



**Automated Malware Detection Process**

Static Analysis → Cluster → Anomaly Detection → Dynamic Analysis → APT Detection

Phase I | Phase II

*Dynamic analysis is leveraged to reduce false positives (from Phase I) and to more accurately identify APTs*

- **Phase I** : Provide us with a computationally efficient way to rapidly identify memory artifacts with anomalous code

- **Phase II**: We leverage call trace information to differentiate between **benign anomalous code** and **malware**

- Security is hard… Why not make it even harder for the adversary?

- Specifically, require the adversary to develop techniques to challenge this memory forensics approach that are both **reliable** (e.g. few bugs) and **portable** (e.g. works across various versions of the OS *and binaries*)

# REFERENCES

1.  Rieck, K., Trinius, P., Willems, C., & Holz, T. (2011). Automatic analysis of malware behavior using machine learning. *Journal of Computer Security*, *19*(4), 639-668.

2.  Stuttgen, Johannes & Cohen, Michael (2013). Anti-forensic Resilient Memory Acquisition. *Digit. Investig., 10*, S105-S115..

3.  Shoshitaishvili, Y., Wang, R., Hauser, C., Kruegel, C., & Vigna, G. (2015, February). Firmalice-Automatic Detection of Authentication Bypass Vulnerabilities in Binary Firmware. In *NDSS*.

4.  Koret, Joxean & Bachaalany, Elias (2015). *The Antivirus Hacker's Handbook.* Wiley Publishing

5.  Carter, K. M., Lippmann, R. P., & Boyer, S. W. (2010, November). Temporally oblivious anomaly detection on large networks using functional peers. In *Proceedings of the 10th ACM SIGCOMM conference on Internet measurement* (pp. 465-471). ACM.

6.  Mosli, R., Li, R., Yuan, B., & Pan, Y. (2016, May). Automated malware detection using artifacts in forensic memory images. In *Technologies for Homeland Security (HST), 2016 IEEE Symposium on* (pp. 1-6). IEEE.

# REFERENCES

6.  Stephens, N., Grosen, J., Salls, C., Dutcher, A., Wang, R., Corbetta, J., ... & Vigna, G. (2016, February). Driller: Augmenting Fuzzing Through Selective Symbolic Execution. In *NDSS* (Vol. 16, pp. 1-16).

7.  Liang, S. C. (2016). Understanding behavioural detection of antivirus.

8.  Anderson, B., Storlie, C., & Lane, T. (2012, October). Improving malware classification: bridging the static/dynamic gap. In *Proceedings of the 5th ACM workshop on Security and artificial intelligence* (pp. 3-14). ACM. Koret, Joxean & Bachaalany, Elias (2015).

9.  Shoshitaishvili, Y., Wang, R., Salls, C., Stephens, N., Polino, M., Dutcher, A., ... & Vigna, G. (2016, May). Sok:(state of) the art of war: Offensive techniques in binary analysis. In *Security and Privacy (SP), 2016 IEEE Symposium on* (pp. 138-157). IEEE.

10. Ravi, C., & Manoharan, R. (2012). Malware detection using windows api sequence and machine learning. *International Journal of Computer Applications*, *43*(17), 12-16.

# Q&A

- How different are equivalent types of memory artifacts ( originating from identical binaries) across multiple hosts?

  - The theory and the empirical results* suggest that memory artifacts are almost identical (+99% similar based on empirical results)

  - *Under the following assumptions

    i. The hosts *are not* under *significant* memory pressure

    ii. Identical versions of the host operating system

# Q&A

- Can we use hashes to determine if memory artifacts across hosts are identical *(e.g. identical explorer.exe process artifacts on Hosts A & B*)?
  - **No.** Chunks of the binary may not be in memory because the OS has likely **paged** those sections to disk in order to efficiently utilize/manage memory
  - Also, the binary is likely **memory mapped**. Therefore, the OS may take a lazy approach by loading a particular chunk when needed
  - As a consequence, the binary on-disk will be different then its image that resides in memory. Memory artifacts across hosts are very likely to also be different.

# APPENDIX

- Machine Learning Primer — A

- Advanced Binary Analysis — B

- Memory Forensics — C

- Binary Vector Generation — D

- Call Trace Vector Generation — E

- Hooking — F

Machine Learning Primer

# APPENDIX A

# MACHINE LEARNING

- Key Concepts

  - **Unsupervised Learning**: Inferring a function to describe hidden structure from "unlabeled" data

  - **Supervised Learning**: Inferring a function from *labeled training data*

  - **Clustering**:  Grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar

  - **Classification**: Identifying to which of a set of categories (sub-populations) a new observation belongs

# MACHINE LEARNING

- Challenges/Steps

  1. *(Non-Trivial)* **Representing** observed/collected data in a meaningful mathematical expressions (e.g. vector)

  2. **Deciding** on a metric for measuring the similarity of observations (e.g. Jaccard Similarity function)

  3. **Selecting** a suitable algorithm that can classify and/or cluster observations appropriately using a supervised or unsupervised approach (e.g. agglomerative hierarchical clustering)

# Machine Learning

- **Example**: Group squares based on similarity of the color types

# MACHINE LEARNING

1.  ## Representing Observations Mathematically

    - We'll represent the contents of each square as a vector where each
      dimension represents the number of occurrences of a color

    $$\begin{pmatrix} \text{Pink} \\ \text{Green} \\ \text{Blue} \end{pmatrix}$$

    - Vector representations



$$\begin{pmatrix} 0 \\ 0 \\ 3 \end{pmatrix} \qquad \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \qquad \begin{pmatrix} 0 \\ 2 \\ 1 \end{pmatrix}$$

2. Metric for measuring similarity of observations

• We'll use the Jaccard Index to measure similarity

$$J(S_i, S_j) = \frac{S_i \cap S_j}{S_i \cup S_j} = \frac{\#\ \text{common color types}}{\#\ \text{total color types}}$$

• Example:



$$J(S_2, S_3) = \frac{S_2 \cap S_3}{S_2 \cup S_3} = \frac{2}{3}$$

3. Selecting a suitable algorithm

- We'll use a hierarchical clustering algorithm

- Depending on the input parameters of the algorithm, the clusters could look like the following

Binary Analysis

# APPENDIX B

# BINARY ANALYSIS

- **Static Analysis**:  Analysis of computer software that is performed without the actual execution of the software code

- **Dynamic Analysis**: Execution of software in an instrumented or monitored manner to garner more concrete information on behavior

# BINARY ANALYSIS

- Static vs. Dynamic Analysis

  - **Static analysis** scales well and can provide better code coverage of a binary

  - **Dynamic analysis** can provide more accurate information on the actual execution behavior of a binary

  - **Static analysis** can produce false execution behavior as code paths may not be reachable during actual execution

  - **Dynamic analysis** can be computationally expensive

# BINARY ANALYSIS

- **Advanced analysis techniques**

  - **Symbolic Execution**: Analysis of a program to determine the necessary inputs needed to reach a particular code path. Variables modeled as symbols

  - **Concolic Execution**: Used in conjunction with symbolic execution to generate concrete inputs (test cases) from symbolic variables to feed into program

  - **Selective Concolic Execution**:  Selectively leverage concolic execution when fuzzing engine gets "stuck" (i.e. unable to generate inputs that can traverse a desired code path)

# BINARY ANALYSIS

- Motivational example for Symbolic Execution

```c
int main(void) {
char buf[32];

char *data = read_string();
unsigned int magic = read_number();

 // difficult check for fuzzing
 if (magic == 0x31337987) {
// Bad stuff
   doBadStuff();
 }
 else if(magic < 100 && magic % 15 == 2 && magic % 11 == 6) {
// Only solution is 17;
   doReallyBadStuff();
 }
 else{
   doBenignStuff();
 }
```

# BINARY ANALYSIS

▪ Motivational example for Symbolic Execution

```c
int main(void) {
char buf[32];

char *data = read_string();
unsigned int magic = read_number();

 // difficult check for fuzzing
 if (magic == 0x31337987) {
 // Bad stuff
   doBadStuff();
 }
 else if(magic < 100 && magic % 15 == 2 && magic % 11 == 6) {
 // Only solution is 17;
   doReallyBadStuff();
 }
 else{
   doBenignStuff();
 }
```

Symbolic execution allows us to figure out the conditions (i.e. magic=0x31337987) to exercise this code path

A more sophisticated code path that can be reached via symbolic execution

# BINARY ANALYSIS

- Analyzing the call traces of example

| Call Trace A |
|---|
| main() |
| doBenignStuff() |

| Call Trace B |
|---|
| main() |
| doReallyBadStuff() |

- **Call Trace A**: Likely result if utilizing traditional emulation techniques to analyze sophisticated malware

- **Call Trace B**: The more useful trace for identifying potential malicious behavior of a binary as a result of applying advanced binary analysis techniques

# APPENDIX C

# MEMORY FORENSICS

- **Defined**: Analysis of a computer's memory dump

- Memory Acquisition
  - Refers to the process of accessing the physical memory
  - *Most critical step in the memory forensics process*
  - Software and hardware tools can be used during acquisition, but we'll focus on the former

- **Rekall** and **Lime** provide open source *acquisition* tools

- **Volatility** and **Rekall** provide open source *analysis* tools

# MEMORY FORENSICS

- Virtual Addressing and Memory Acquisition

**Process A**

2148ABCD  0x12345

0x12467

**Process B**

DEADBEAF  0x325410

2148ABCD  0x57890

**Physical Memory**

0x1FF9612

DEADBEAF

2148ABCD  0x42FF9612

FFFFFFFF  0x62FF9612

**Disk**

# MEMORY FORENSICS

- Virtual Addressing and Memory Acquisition (cont'd)

**Process A**

**Physical Memory**

0x1FF9612

**Disk**

0x12467

- Physical space may be smaller than virtual address space.

- Less recently used memory blocks (a.k.a. pages) are moved to disk

- **Important**: *Only data that is in physical memory during acquisition can be acquired; **paged data is unavailable***

# Memory Forensics

- **Anti-Forensics**

  - Any attempt to compromise the availability or usefulness of evidence to the forensic process

  - Techniques include

    i. **Substitution Attack**: Data fabricated by the attacker is substituted in place of valid data during the acquisition

    ii. **Disruption Attack**: Disrupt the acquisition process

  - Proof of Concepts:

    i. "ShadowWalker" @ Blackhat 2005

    ii. "Low Down and  Dirty"  @ Blackhat 2006

    iii. "Defeating Windows  Forensics" @Fahrplan 2012

  - *The presented approach is resilient to anti-forensics techniques due to information asymmetry on the side of the defender*

# Generating a Vector Representation of a Binary

# APPENDIX D

# BINARY VECTOR GENERATION



- Binary Vector Generation
  i. **Load** the dumped binary into a disassembler (e.g. IDA or Binary Ninja) and extract the executable portion of binary
  ii. **Lift** the executable portion of binary to a RISC-like intermediate representation (ir)
  iii. **Create** a binary behavior report that categorizes each ir statement
  iv. **Generate** a set of fixed-sized linear sequences (w.r.t. address space) of ir statements that will be referred to as behavior sequences
  v. **Create** a map that maps each unique sequence to a row in a vector

# BINARY VECTOR GENERATION

| Disassembler | → | RISC-like Instructions | → | Behavior Report |

| **main()** |
| --- |
| push ebp |
| mov rbp, rsp |
| sub rsp, 0x10 |

| **RISC-like (vex ir )** |
| --- |
| t0 = GET:I64(bp) |
| t3 = GET:I64(rsp) |
| t2 = Sub64(t3, 0x008) |
| PUT(rsp) = t2 |
| St1e(t1) = t0 |
| t3 = GET:I64(rsp) |
| PUT(bp) = t3 |
| t4 = GET:I64(rsp) |
| t5 = 0x10 |
| t6 = Sub64(t4,t5) |
| PUT(rsp) = t6 |

| **Report** |
| --- |
| reg_access |
| reg_access |
| arithmetic |
| reg_access |
| store |
| reg_access |
| reg_access |
| reg_access |
| other |
| arithmetic |
| reg_access |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| **reg_access** |
| **reg_access** |
| **arithmetic** |
| reg_access |
| store |
| reg_access |
| reg_access |
| reg_access |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| reg_access |
| **reg_access** |
| **arithmetic** |
| **reg_access** |
| store |
| reg_access |
| reg_access |
| reg_access |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

| SEQ 2 |
|---|
| reg_access |
| arithmetic |
| reg_access |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| reg_access |
| reg_access |
| **arithmetic** |
| **reg_access** |
| **store** |
| reg_access |
| reg_access |
| reg_access |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

| SEQ 2 |
|---|
| reg_access |
| arithmetic |
| reg_access |

| SEQ 3 |
|---|
| arithmetic |
| reg_access |
| store |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| reg_access |
| reg_access |
| arithmetic |
| **reg_access** |
| **store** |
| **reg_access** |
| reg_access |
| reg_access |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

| SEQ 2 |
|---|
| reg_access |
| arithmetic |
| reg_access |

| SEQ 3 |
|---|
| arithmetic |
| reg_access |
| store |

| SEQ 4 |
|---|
| reg_access |
| store |
| reg_access |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| reg_access |
| reg_access |
| arithmetic |
| reg_access |
| **store** |
| **reg_access** |
| **reg_access** |
| reg_access |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

| SEQ 2 |
|---|
| reg_access |
| arithmetic |
| reg_access |

| SEQ 3 |
|---|
| arithmetic |
| reg_access |
| store |

| SEQ 4 |
|---|
| reg_access |
| store |
| reg_access |

| SEQ 5 |
|---|
| store |
| reg_access |
| reg_access |

# BINARY VECTOR GENERATION

| Behavior Report | → | Behavior Sequences |
|---|---|---|

| Report |
|---|
| reg_access |
| reg_access |
| arithmetic |
| reg_access |
| store |
| **reg_access** |
| **reg_access** |
| **reg_access** |
| other |
| arithmetic |
| reg_access |

| SEQ 1 |
|---|
| reg_access |
| reg_access |
| arithmetic |

| SEQ 2 |
|---|
| reg_access |
| arithmetic |
| reg_access |

| SEQ 3 |
|---|
| arithmetic |
| reg_access |
| store |

| SEQ 4 |
|---|
| reg_access |
| store |
| reg_access |

| SEQ 5 |
|---|
| store |
| reg_access |
| reg_access |

| SEQ 6 |
|---|
| reg_access |
| reg_access |
| reg_access |

- Total possible behavior sequence permutations: ~ $2^{77}$

  - Number of instruction categories (e.g. store and branch): 14

  - Length of behavior sequence: 20

  - $14^{20}$ ~ $2^{77}$

- **Naive Approach**

  - Create a $14^{20}$ dimensional vector to express binary behavior

  - ***Each vector dimension maps to a unique behavior sequence***

  - Number stored at dimension $k$ is the number of times behavior sequence occurs in executable

# BINARY VECTOR GENERATION

- **Naive Approach (continued...)**



**Index 0**

reg_access
reg_access
reg_access
reg_access
.......
reg_access

20 category items

**Index k**

xxx
xxx
xxx
xxx
......... .
xxx

Number of occurrences *m* of behavior sequence *k*

0    **3**
1    **6**

*k*    **m**

$14^{20}$-1    **0**

**Behavior Vector**

**Index $14^{20}$ -1**

call
call
call
call
......... .
call

- **Naive Approach (continued...)**

  - Not very practical to implement directly

  - Fortunately, we can do better

- **Key observation**: Vector is sparse in that most of the dimensions will store the number '0'

- **Better Approach**

  - Only store the non-zero elements in memory

  - So if a binary has N total ir statements, then we only need to store at most N-`window_length` elements in memory

Efficient Diffing Algorithm

# APPENDIX D

# EFFICIENT DIFFING ALGOIRHTM



Step 2: Compute the similarity of a pair of vectors

explorerA.exe    explorerB.exe

explorerA.exe → explorerB.exe → **Similarity Function** → .867

*Similarity function takes as input two vectors and produces a value between 0 and 1*

- Selecting a similarity function

  - For convenience, we'll use the Jaccard Index

  - The Jaccard Index has the following interpretation:

$$J(S_i, S_j) = \frac{S_i \cap S_j}{S_i \cup S_j} = \frac{\#\text{common behavior sequences}}{\#\text{total behaviors sequences}}$$

**(See Appendix A for a simple example using the Jaccard Index)**

Call Trace Vector Generation

# APPENDIX E

# CALL TRACE VECTOR GENERATION

Disassembler → RISC-like instructions → Call Trace Generation (Dyn. Analysis) → Call Trace Vector

- Call Trace Vector Generation
  i. **Load** the dumped binary into a disassembler (e.g. IDA or Binary Ninja) and extract all sections of the binary
  ii. **Lift** the executable portion of binary to a RISC-like intermediate representation (ir)
  iii. **Perform** Dynamic analysis on the ir to generate call traces
  iv. **Generate** a call trace vector that maps unique call trace sequences into a row of the vector

# CALL TRACE VECTOR GENERATION

Host F

### Call Trace (Code Path 1)

```
main()
LoadLibraryW()
RegisterShellHook()
GetTokenInformation()
LsaOpenPolicy()
VirtualAlloc()
CreateEventW()
LogoffWindowsDialog()
```

### Call Trace (Code Path 2)

```
main()
LoadLibraryW()
RegisterShellHook()
GetTokenInformation()
LsaOpenPolicy()
VirtualAlloc()
CreateEventW()
LogoffWindowsDialog()
```

• • •

### Call Trace (Code Path N)

```
main()
LoadLibraryW()
RegisterShellHook()
GetTokenInformation()
LsaOpenPolicy()
VirtualAlloc()
CreateEventW()
LogoffWindowsDialog()
```

# CALL TRACE VECTOR GENERATION

Call Trace() → Trace Sequences

| Call Trace (Code Path n) |
|---|
| **main()** |
| **LoadLibraryW()** |
| **RegisterShellHook()** |
| GetTokenInformation() |
| LsaOpenPolicy() |
| VirtualAlloc() |
| CreateEventW() |
| LogoffWindowsDialog() |

| Trace Seq 1 |
|---|
| main() |
| LoadLibraryW() |
| RegisterShellHook() |

# CALL TRACE VECTOR GENERATION

Call Trace() → Trace Sequences

## Call Trace (Code Path n)

main()

**LoadLibraryW()**

**RegisterShellHook()**

**GetTokenInformation()**

LsaOpenPolicy()

VirtualAlloc()

CreateEventW()

LogoffWindowsDialog()

## Trace Seq 1

main()

LoadLibraryW()

RegisterShellHook()

## Trace Seq 2

LoadLibraryW()

RegisterShellHook()

GetTokenInformation()

# CALL TRACE VECTOR GENERATION

| Call Trace() | → | Trace Sequences |
|---|---|---|

| Call Trace (Code Path n) |
|---|
| main() |
| LoadLibraryW() |
| **RegisterShellHook()** |
| **GetTokenInformation()** |
| **LsaOpenPolicy()** |
| VirtualAlloc() |
| CreateEventW() |
| LogoffWindowsDialog() |

| Trace Seq 1 |
|---|
| main() |
| LoadLibraryW() |
| RegisterShellHook() |

| Trace Seq 2 |
|---|
| LoadLibraryW() |
| RegisterShellHook() |
| GetTokenInformation() |

| Trace Seq 3 |
|---|
| RegisterShellHook() |
| GetTokenInformation() |
| LsaOpenPolicy() |

# CALL TRACE VECTOR GENERATION

Call Trace() ➡ Trace Sequences

## Call Trace (Code Path n)

main()

LoadLibraryW()

RegisterShellHook()

**GetTokenInformation()**

**LsaOpenPolicy()**

**VirtualAlloc()**

CreateEventW()

LogoffWindowsDialog()

## Trace Seq 1

main()

LoadLibraryW()

RegisterShellHook()

## Trace Seq 2

LoadLibraryW()

RegisterShellHook()

GetTokenInformation()

## Trace Seq 3

RegisterShellHook()

GetTokenInformation()

LsaOpenPolicy()

## Trace Seq 4

GetTokenInformation()

LsaOpenPolicy()

VirtualAlloc()

- **Trace Call Vector**



Index 0

```
main()
LoadLibraryW()
RegisterShellHook()
```
....

Index $k$

xxx
xxx
xxx
xxx
......... .
xxx

Number of occurrences $m$ of trace sequence $k$

Index $14^{20}-1$

```
call
call
call
call
......... .
call
```

Trace Call Vector

# CALL TRACE VECTOR GENERATION

**Trusted**

**Call Behavior Vector**

**Unknown**

Host F

Hooking

# APPENDIX F

# HOOKING

**CreateFileA (Original)**

```
mov  edi, edi
push ebp
mov ebp, esp
```

Modified Instruction

**CreateFileA (Hooked)**

```
jmp MonitorCode
push ebp
mov ebp, esp
```

**Monitor Code**

```
inc dword ptr [edx+8]
.......
jmp CreateFileA + 5
```
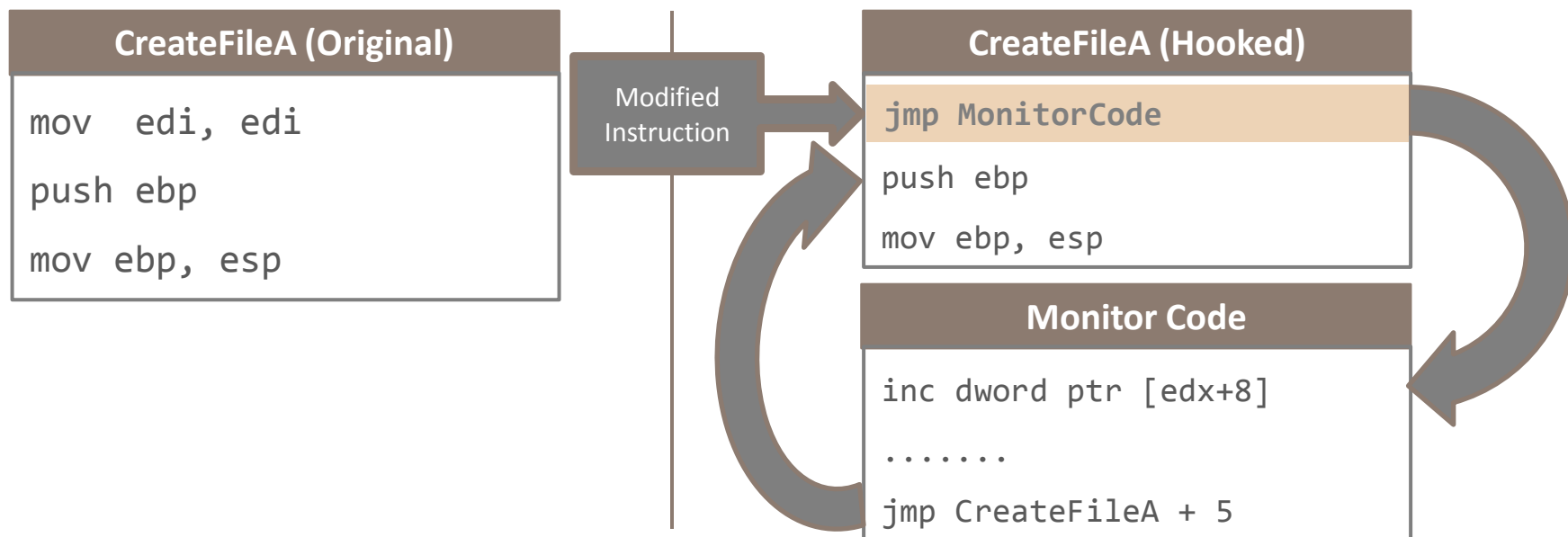
▪ **Monitoring Behavior w/ Hooks**

- Detouring a number of common APIs (e.g. `CreateFile`) to ensure monitoring code is executed before actual code

- Depending on a set of rules ( typically dynamic), the ***monitor code blocks, allows, or reports execution of API***