

Intro: Malware behavior variability

Example: Ramnit worm

Exploits CVE-2013-3660

Only on vulnerable Windows 7, not running in admin mode

Makes ~100 mutexes creation calls

Dataset

Passive recording of:
26K suspicious programs

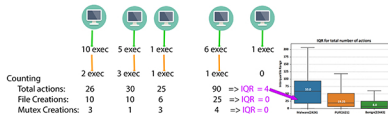
7.6M execution traces
5.6M real users' machines
113 countries

All the year of **2018**

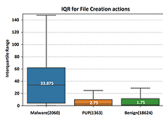
Action Type	Filename	File path	Sample Hash	Execution Date	Machine ID
File Create	setup.exe	C:\MSI_PROFILES	AAA	1	abc
Mutex Create	mta12345	-	ABC	5231021	abd

(Split the data by sample hash)

Action Type	Filename	File path	Sample Hash	Execution Date	Machine ID
File Create	setup.exe	C:\MSI_PROFILES	AAA	1	abc
Mutex Create	mta12345	-	ABC	5231021	abd



>50% of the **malware** have **59** missing/additional actions across all the machines in the **first week** of their appearance



File Creation is the **major source** of machine-induced variability.

>50% of the malware have **33** missing/additional file creations

		Median			75 th percentile		
		Mal	PUP	Ben	Mal	PUP	Ben
File	Path	4	1	-	10	3	2
	Name	25	2	1	49	8	8
	Ext.	3	1	-	5	2	1

Variability is in the **number of unique file names** not paths. Malware write in the same number of unique file paths.

Case study: Glupteba

	Jaccard index	IQR
File name	0	0
Mutex name	0.2	2

Mutex values appeared in:

h48yorbq6rm87zot — all the machines

ZoneCacheCounterMutex — in ~50%

ZoneAttributeCacheCounterMutex — in ~50%

2 random values appear in only 1 machine

$$\text{Jaccard} = \frac{\text{number of values appearing in all machines}}{\text{number of all unique values}} = \frac{1}{5} = 0.2$$

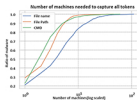
$$\text{IQR} = (\text{num. unique mtx. names in 75 percentile}) - (\text{those in 25 percentile}) = 3 - 1 = 2$$

Mutex name is a better candidate for building signatures for this specific malware sample

Invariant behavior analysis

```
logsource:
  category: process_creation
  product: windows
detection:
  selection:
    CommandLine: "mon! es bypass"
condition: selection
```

Measure prevalence of tokens in parameter values for 3 most used parameters.



Q1: How to capture all malware tokens?

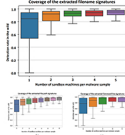
File name tokens are more challenging to capture than the file path and CMD tokens.

Q2: How to maximize detection with minimum number of sandboxes?

File name tokens from **3** machines reach **high detection**, more machines provide diminishing ret.

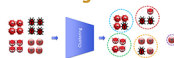
File path tokens need to be extracted from at least **7** machines

CMD line also requires 3 to 4 machines at most



Variability VS malware solutions

Clustering

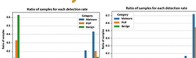


Intuition: All sample's executions will be in the same cluster.

From all malware samples

- 67% fall in 1 cluster
- 27% fall in 2 clusters
- 5% fall in 3 clusters
- 1% fall in 4 different clusters

Anomaly detection



Using all benign exec. Using 1 benign exec.

Using all benign executions we get **lower** FP but also **lower** detection.

- When employing malware analysis methods:
 - Use multiple executions of the same malware samples
 - Use multiple executions of the benign samples

Measuring cross-machine variability

Action Type	Filename	File path	Sample Hash	Execution Date	Machine ID
File Create	setup.exe	C:\MSI_PROFILES	AAA	1	abc
Mutex Create	mta12345	-	ABC	5231021	abd

