



# AndroOBFS: Time-tagged Obfuscated Android Malware Dataset with Family Information

---

Saurabh Kumar<sup>1</sup>, Debadatta Mishra<sup>1</sup>, Biswabandan Panda<sup>2</sup>, and Sandeep K. Shukla<sup>1</sup>

<sup>1</sup>Indian Institute of Technology Kanpur

<sup>2</sup>Indian Institute of Technology Bombay

**MSR 2022**

# Motivation

- ❑ Malware analysis system require labeled dataset
  - Non-obfuscated (Drebin, AMD, RmvDroid, etc..)
  - Obfuscated (PRAGuard)
  
- ❑ PRAGuard dataset
  - Outdated samples till March 2013
  - Does not contains family information
  
- ❑ Needed a new obfuscated malware dataset
  - Tagged with time
  - With family information

# AndroOBFS Dataset

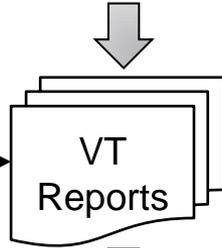
- ❑ Samples collected from
  - AndroZoo and VirusShare
  - For year 2018, 2019 and 2020
- ❑ Tools used
  - AVClass: Labeling with family
  - Obfuscapk: To obfuscate
- ❑ 16279 obfuscated samples in six categories
  - 14579 Familial
  - Spread across 158 families

Sources	Year	#Samples	
		Obfuscated	Obfuscated (#Family <sub>≥2</sub> )
AndroZoo + VirusShare	2018	6525	5794 (136)
	2019	8313	7505 (94)
AndroZoo	2020	1441	1280 (44)
<b>Total</b>		<b>16279</b>	<b>14579 (158)</b>

# Dataset Creation Process

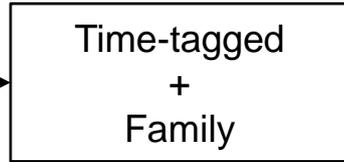
**VirusShare + AndroZoo**

 **VIRUSTOTAL**

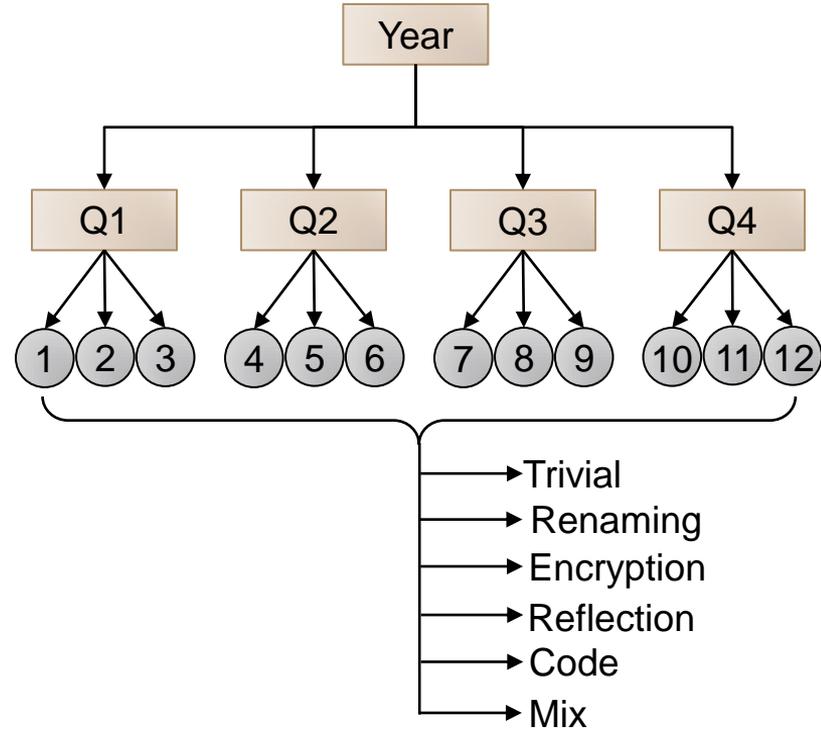
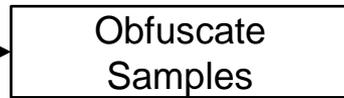


$\forall VT^+ \geq 10$

**Dex Date**  
+  
**AVClass**

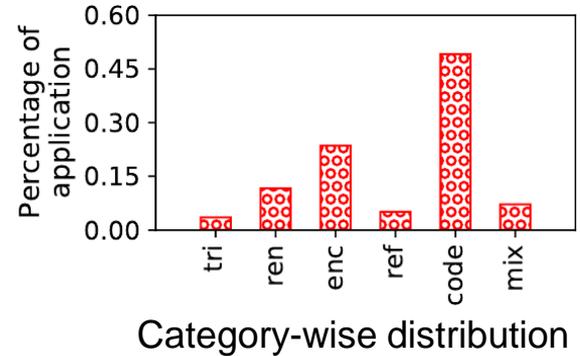


**Obfuscapk**



# Malware Distribution

- ❑ 16279 samples distributed among six categories
  - ~50% fall under code obfuscation
  - ~23% in encryption



# Usage Scenarios

- ❑ Designing robust malware detection and classification system
- ❑ Studying the efficacy of existing analysis systems
- ❑ Temporal/longitudinal study of an analysis system

# Conclusion

16279 obfuscated samples in six categories

14579 familial obfuscated samples distributed among 158 families

Tagged with time

Spanning over three years from 2018 to 2020

Thank You