

Information Security Centre of Excellence

Automated Semi-supervised Authorship Attribution of Android Binaries Hugo Gonzalez, Natalia Stakhanova, Ali A. Ghorbani Faculty of Computer Science, University of New Brunswick



Computer Science

The Problem

- Malware on the Android platform has increased exponentially. Growing more than 600% in the last two years.
- New threats appear each day in the form of botnet, ransomware, adware or even worse: all together
- Typical methods to perform malware analysis focus on binaries it self, creating signatures or behavioral profiles.



Proposed solution

- \succ It is our believe that only a reduced set of skilled authors produce primary strands of the malware, meanwhile the available samples are repackaged, reused, recompiled or evolved versions from the main strand.
- Our proposed approach looks for similarities in Android binary code from authors perspective.
- It involves 3 step process:
- Creation and evaluation of benchmark profiles
- Creation and evaluation of incremental profiles
- Emergent behavior layer, where a large number of apps are analyzed.



Contributions

Android Authors Dataset

- Method to classify and cluster on top of Random Forest algorithm
- Approach to create label profiles
- Large scale evaluation

Future work

Name Authors dataset

Adware d

- Evaluate complementary features
- Improve random forest performance using distributed task
- Use small number of files to create possible profiles.
- Web service to analyze apps.

Datasets		Feature Extraction		
Name Authors dataset	Table 3 Apps 1,428	3: Datasets used Description Collected from eight different markets this dataset contains in	Values extracted from apps	Semi -Supervised learning
Fdroid dataset Drebin partial dataset	1,395 3,181	formation about 33 known au- thors. Open source apps without adver- tisement libraries. From the original dataset, we keep only families with more then 20 apps, to explore relations be-	(classes, metods) +	Random Forest
Adware dataset Ransomware dataset GooglePlay- 2015 dataset	211 136 4,574	 tween families and authors. Three families and authors. Three families related between them. They trojanize legit apps. Ransomware software usually contain specialized behavior. It includes four families. Apps collected from Google Play in middle 2015 from top popular 	Array related opcodes	3gram, 32 768 features, 10 files per groups and 8 bins.

Total Unseen apps	33,153		
Attributed apps	4,830	14.57%	
Correctly attributed apps *	3,371	69.79%	10.17%
Unattributed apps	28,323	85.43%	
Related without label	11,839	35.71%	
Non related apps	16,484	49.72%	
New created label profile	1,147		
With more than 200 apps	3	Plankton	
With more than 100 but less than 200	12	Googleplay	

* Apps are considered correctly attributed when they belong to the same malware family or from the same market. Further investigation is needed.