

ISCX

Information Security  
Centre of Excellence

# DroidKin: Lightweight Detection of Android Apps Similarity

Hugo Gonzalez, Natalia Stakhanova, Ali A. Ghorbani  
Faculty of Computer Science, University of New Brunswick



## Our work

DroidKin, a robust approach for detection of Android apps similarity. Based on a set of characteristics derived from binary and meta data accompanying apk files, DroidKin is able to detect similarity among applications under various levels of obfuscation. DroidKin performs analysis pinpointing similarities between applications and identifying their relationships.

## Contributions

- Develop a lightweight and robust approach to app similarity detection
- Investigate the impact of lightweight features: meta-data, meta-information and content, on detection of plagiarized apps

## Features

- Meta-data** that accompanies each .apk file. The following features were extracted: serial number, developer's name, size of an .apk file. This descriptive information serves as a reference for analyzed applications.
- Meta-information** that characterizes .apk file contents. These features include timestamps of the content files; number of internal .apk, .zip, .java, .jar, images, libraries and binary files found within a container .apk file; size of .dex file; hash value of .dex file (md5); number of files in .apk file; list of their names with corresponding hash values (md5); number of unique timestamps from the files.
- N-grams** characterizing the .dex file, using opcodes only, opcodes and operands and bytecode.

## Datasets

- Created a dataset with different levels and types of malware app obfuscation (720 apk files)



- Android Malware Genome Project (1260 apps).

**DREBIN**

- dataset. (5560 apps).

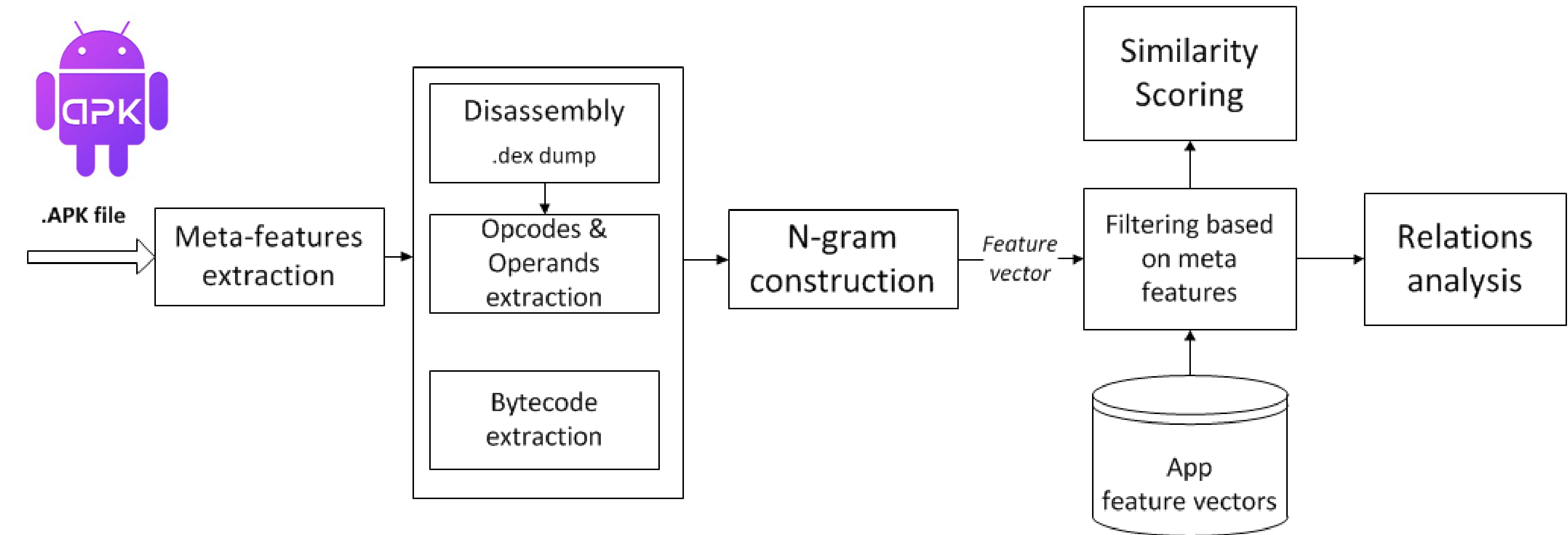
## Scoring

The pair-wise similarity between related apps is calculated using a variation of a similarity measure Simplified Profile Intersection (SPI). The metric has been proposed for evaluation of source code author profiles based on frequency analysis. Given Simplified Profiles  $SP_i$  and  $SP_j$ , the similarity distance, called Similarity Profile Intersection SPI is the size of intersection between these profiles. We use meta-data, meta-information and n-grams to calculate this.

$$SPI_n = |SP_i \cap SP_j|$$

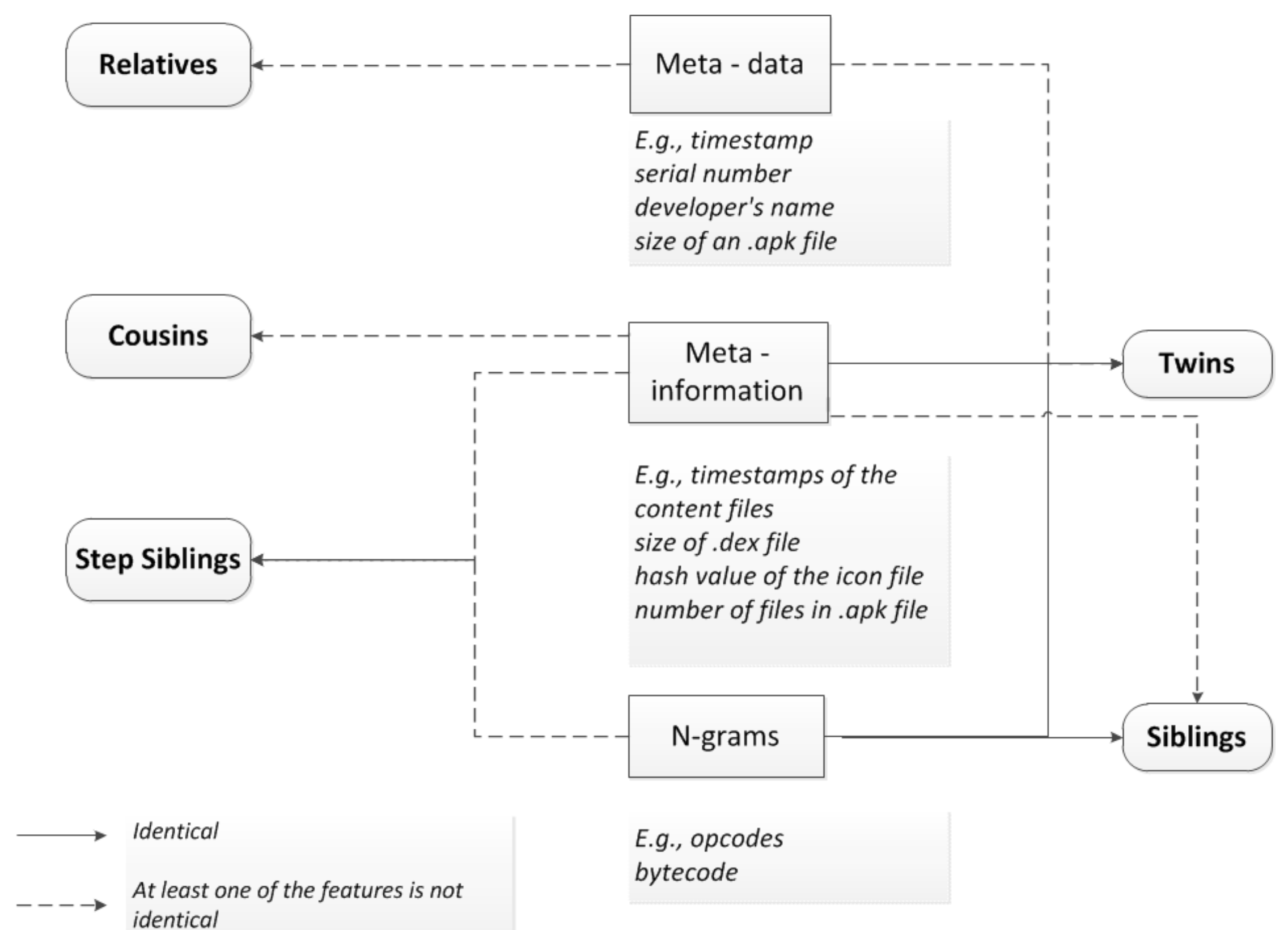
- SPI<sub>n</sub>** = SPI for the frequency of the ngrams extracted from dex file.
- SPI<sub>c</sub>** =  $.6 * SPI$  for the content of the files +  $.4 * SPI$  for the names of the files.

## System Design



## Relations

- Twin:**  $SPI_n = 100\%$  &  $SPI > 95\%$
- Sibling:**  $SPI_n = 100\%$  &  $SPI > 60\%$
- Step-sibling:**  $SPI_n > 60\%$  &  $SPI > 60\%$
- Cousin:**  $SPI_n < 60\%$  &  $SPI > 60\%$
- False step-sibling:**  $SPI_n > 60\%$  &  $SPI < 60\%$



## Results

