

# FINGERPRINTING ANDROID PACKAGING GENERATING DNA'S FOR MALWARE DETECTION

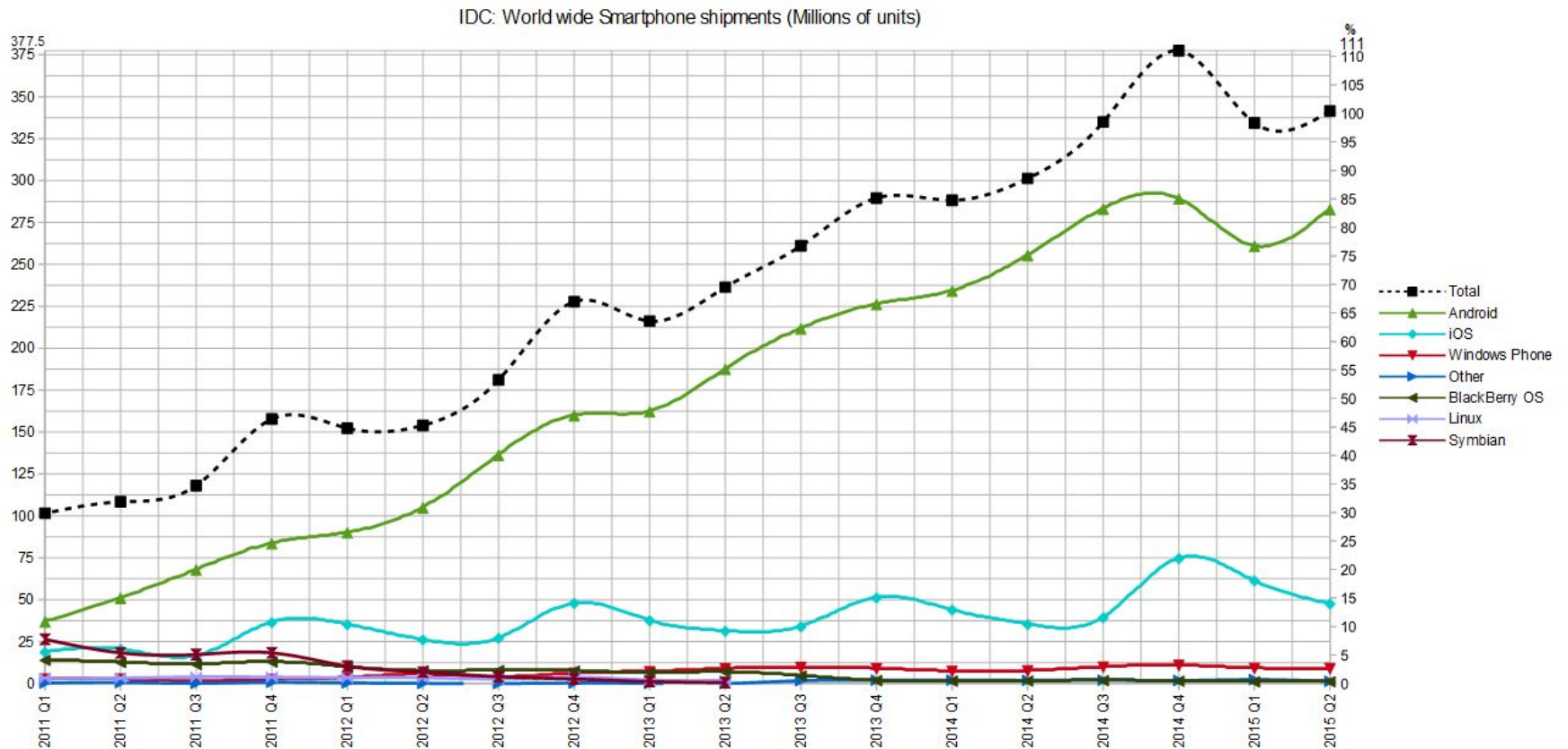
**ElMouatez B. Karbab**, Mourad Debbabi, Djedjiga Mouheb

August 8, 2016

DFRWS USA Conference 2016

# Motivation

2

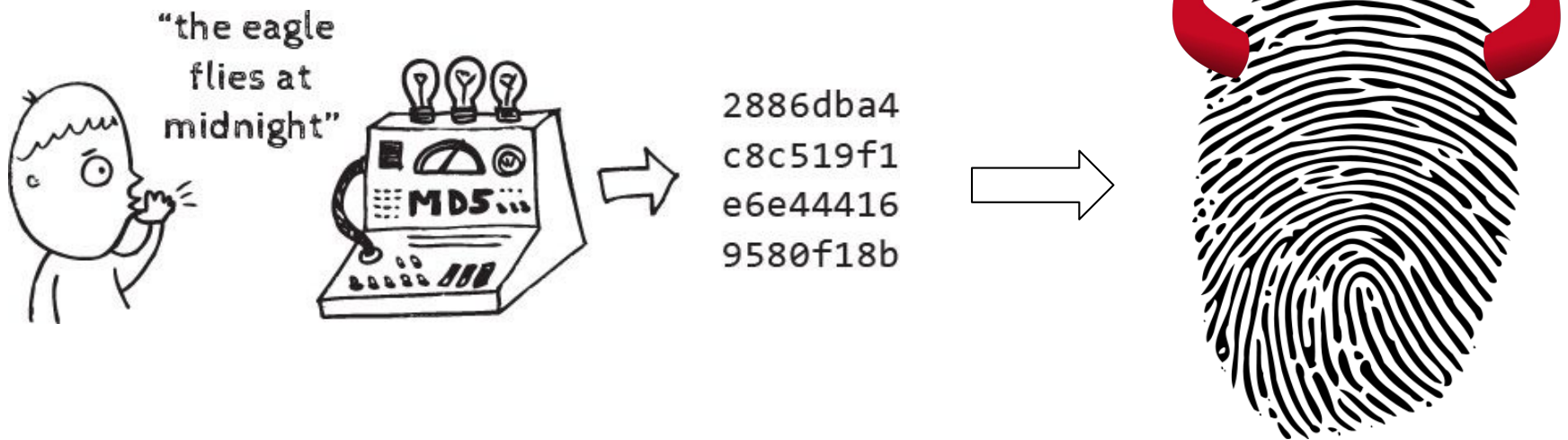


[https://commons.wikimedia.org/wiki/File:MobileOS\\_market\\_share\\_till\\_2014\\_Q2.png](https://commons.wikimedia.org/wiki/File:MobileOS_market_share_till_2014_Q2.png)

# Motivation (cont'd)

3

## □ Crypto Hash: MD5, SHA1 ...



Malware Fingerprint  
Using Cryptographic  
Hashing

# Motivation (cont'd)

4

- Fuzzy Hash: ssdeep, mv-Hash ...



SHA256: 7dcb02d16fe53b3de5536428a4bd3436e41e076760727909eca3c4b1c97985a0

File name: 7dcb02d16fe53b3de5536428a4bd3436e41e076760727909eca3c4b1c97985a0.apk

Detection ratio: 42 / 55

Analysis date: 2016-07-14 07:07:56 UTC ( 3 weeks, 3 days ago )

Analysis File detail Additional information Comments 0 Votes Behavioural information

## File identification

MD5	7c8877f8a6afc41dbacfa3db4931da1
SHA1	fa9d92d1aad565c7ef24222bc9ea8de5c426decf
SHA256	7dcb02d16fe53b3de5536428a4bd3436e41e076760727909eca3c4b1c97985a0
ssdeep	98304:tZHmSJaxGUtvFg94QzFezkJUz8c431vgquHQ:tZHmWgGUtdg9zzQzIU3+dd
File size	3.4 MB ( 3583428 bytes )
File type	Android
Magic literal	Zip archive data, at least v2.0 to extract
TrID	Android Package (73.9%) Java Archive (20.4%) ZIP compressed archive (5.6%)
Tags	apk android contains-elf

Malware Fingerprint  
Using Fyzzing  
Hashing

# Current Fuzzy Hashing Problem

5

- ❑ Ignoring the underneath structure and semantics of the malicious package.
- ❑ Fuzzy hashing suffers from its single fingerprint bounded with a maximum size.
- ❑ Android Packaging Noise.

# Objectives

6

- Develop a more accurate, yet broad, fuzzy fingerprinting technique for Android malicious apps packages.
- Design and implement a framework for Android malware detection and family attribution on top of the developed fuzzy fingerprint.

# APK-DNA

7

- Fuzzy fingerprint covers the underneath Android app structure.
- An automatic framework for Android malware detection using APK-DNA
  - Peer-matching approach
  - Family-fingerprinting approach

# Android Package (APK)

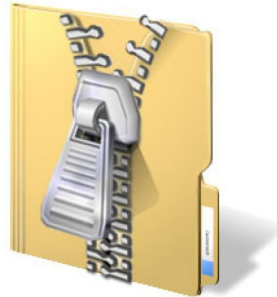
8

## Contents

Android  
Package (APK)



Unzipping

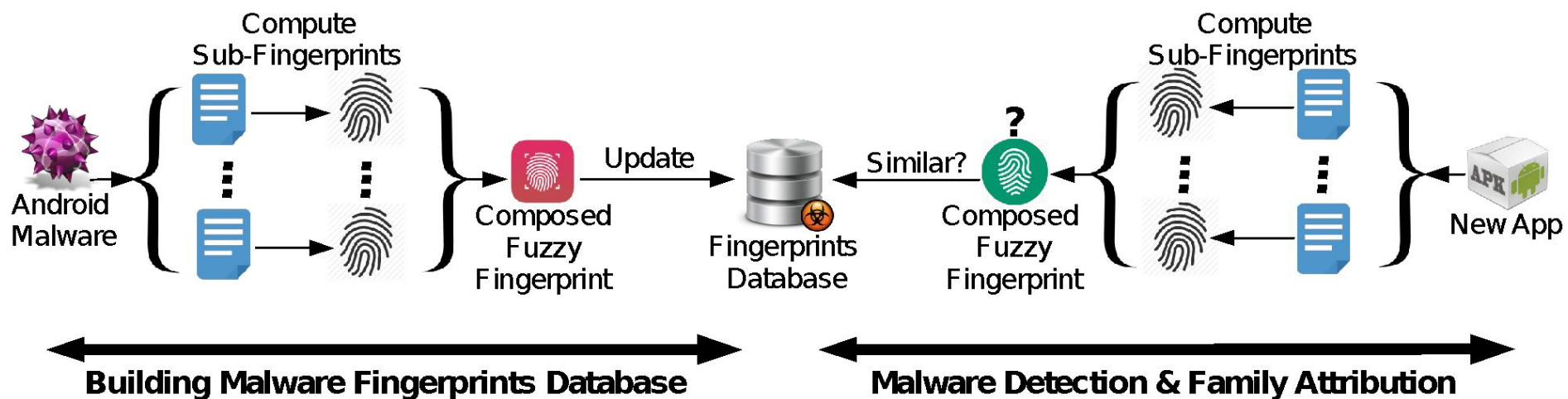


Name	Size	Compressed
META-INF	3 Files	
ALIAS_NA.RSA	1.3 KiB	1.0 KiB
ALIAS_NA.SF	1.7 KiB	839 B
MANIFEST.MF	1.6 KiB	733 B
assets	1 File	
libparser.so	291.7 KiB	291.7 KiB
lib	1 Folder	
armeabi	3 Files	
res	13 Folders	
drawable-hdpi	1 File	
drawable-ldpi	1 File	
drawable-mdpi	1 File	
layout	1 File	
raw	1 File	
raw-de	1 File	
raw-es	1 File	
raw-fr	1 File	
raw-it	1 File	
raw-ja	1 File	
raw-zh-CN	1 File	
raw-zh-TW	1 File	
xml	1 File	
AndroidManifest.xml	4.1 KiB	1.2 KiB
classes.dex	146.5 KiB	67.8 KiB
resources.arsc	18.2 KiB	5.2 KiB



# Approach

9



# Android Package (APK)

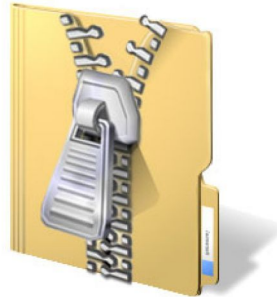
10

## Contents

Android  
Package (APK)



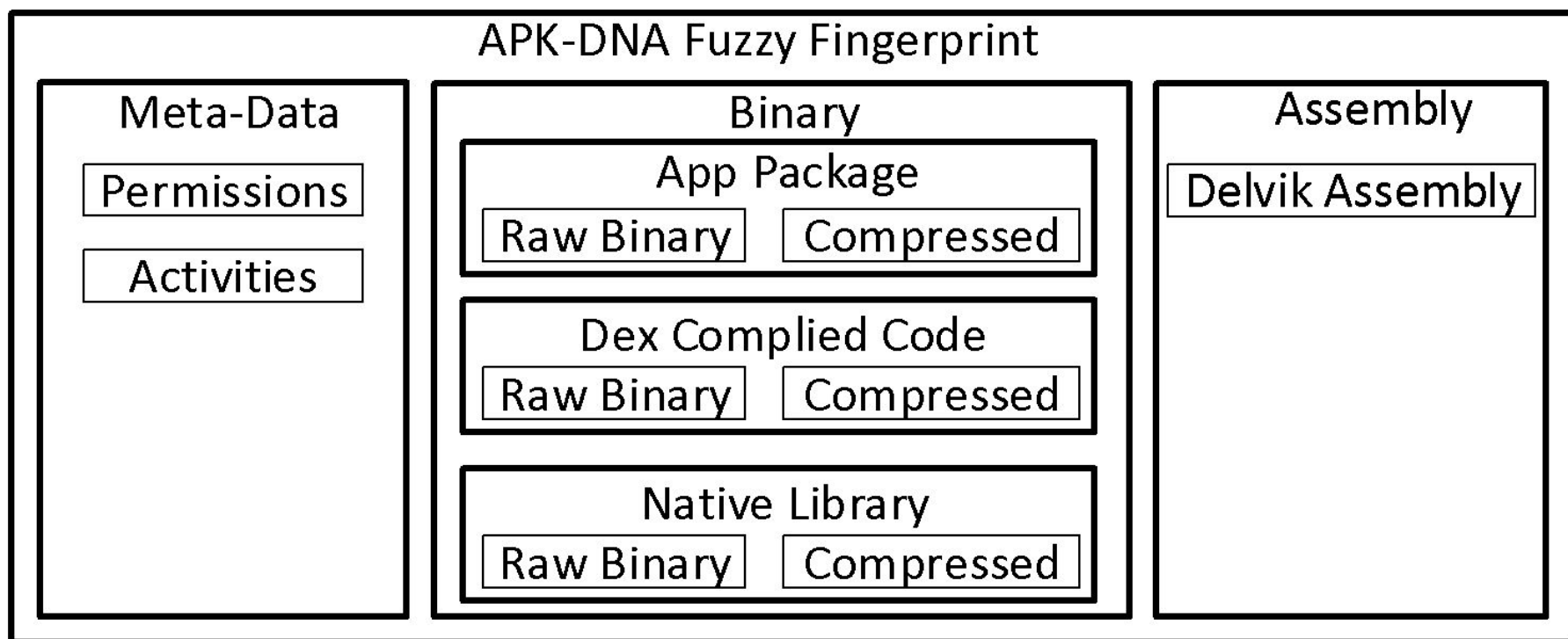
Unzipping



Name	Size	Compressed
META-INF	3 Files	
ALIAS_NA.RSA	1.3 KiB	1.0 KiB
ALIAS_NA.SF	1.7 KiB	839 B
MANIFEST.MF	1.6 KiB	733 B
assets	1 File	
libparser.so	291.7 KiB	291.7 KiB
lib	1 Folder	
armeabi	3 Files	
res	13 Folders	
drawable-hdpi	1 File	
drawable-ldpi	1 File	
drawable-mdpi	1 File	
layout	1 File	
raw	1 File	
raw-de	1 File	
raw-es	1 File	
raw-fr	1 File	
raw-it	1 File	
raw-ja	1 File	
raw-zh-CN	1 File	
raw-zh-TW	1 File	
xml	1 File	
AndroidManifest.xml	4.1 KiB	1.2 KiB
classes.dex	146.5 KiB	67.8 KiB
resources.arsc	18.2 KiB	5.2 KiB

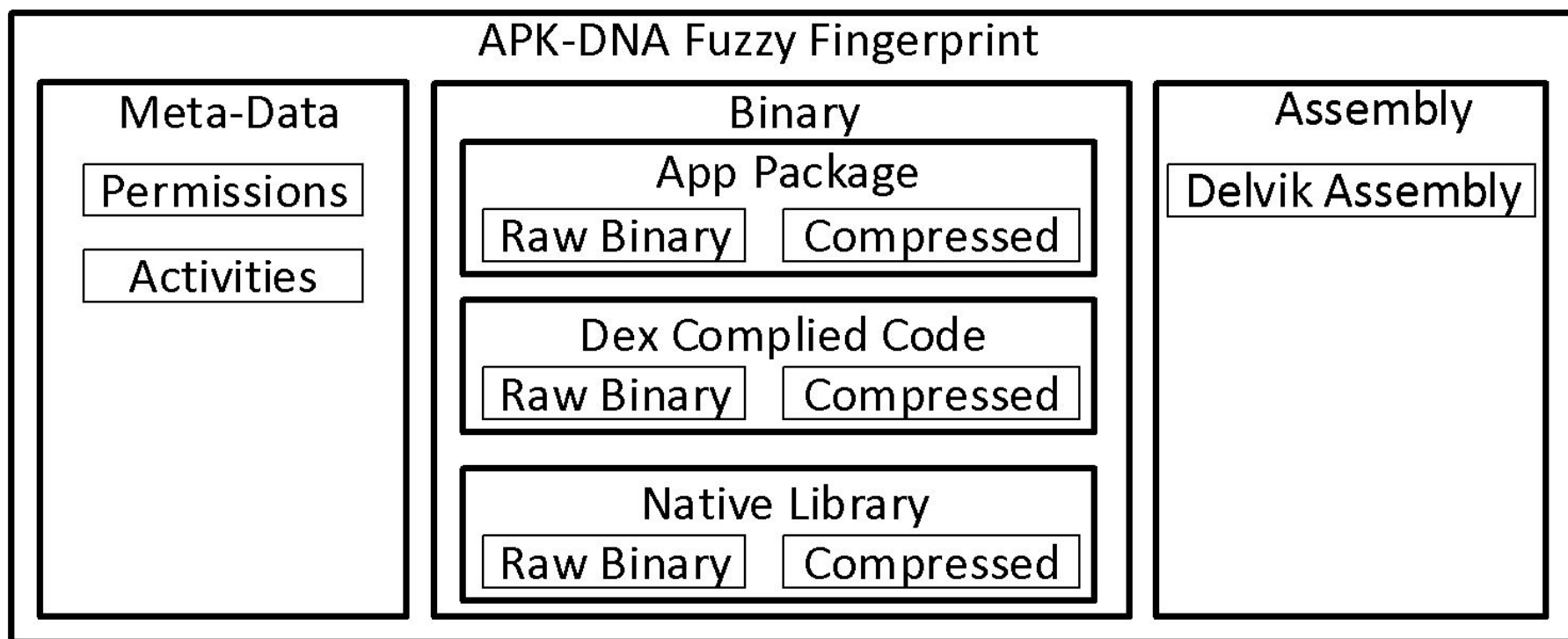
# APK-DNA Structure

11



# APK-DNA Structure

12



# APK-DNA Structure (cont'd)

13

return-void const/4 sput-object return-void invoke-  
direct return-void invoke-static move-result return  
sget-object if-nez new-instance invoke-direct sput-  
object sget-object move-object move-object move-  
object move-object move-object invoke-virtual/range

→ Bigrams

3-grams

5-grams

64 65 78 0a 30 33 35 00 c9 d5 ca 3f 19 a1 31 ee 09 b3  
88 70 ff 1c cf cc fe f9 43 c8 12 34 02 4d 6c 56 02 00 70  
00 00 00 78 56 34 12 00 00 00 00 00 00 00 00 9c 55 02

First Instructions and Bytes of AnserverBot Malware

# Feature Hashing Technique

14

---

## Algorithm 1: Feature Vector Computation

---

```
input  : N -grams: Set,  
        L : Feature Vector Length  
output: Binary Feature Vector  
features_vector = new bitvector[L];  
for Item in N-grams do  
    | H = hash(Item) ;  
    | feature_index = H mod L ;  
    | features_vector[feature_index] = 1 ;  
end
```

---

# Sub-Fingerprint Similarity

15

$$Jaccard(X, Y) = \frac{|X \cap Y|}{|X \cup Y|}$$
$$0 \leq Jaccard(X, Y) \leq 1$$

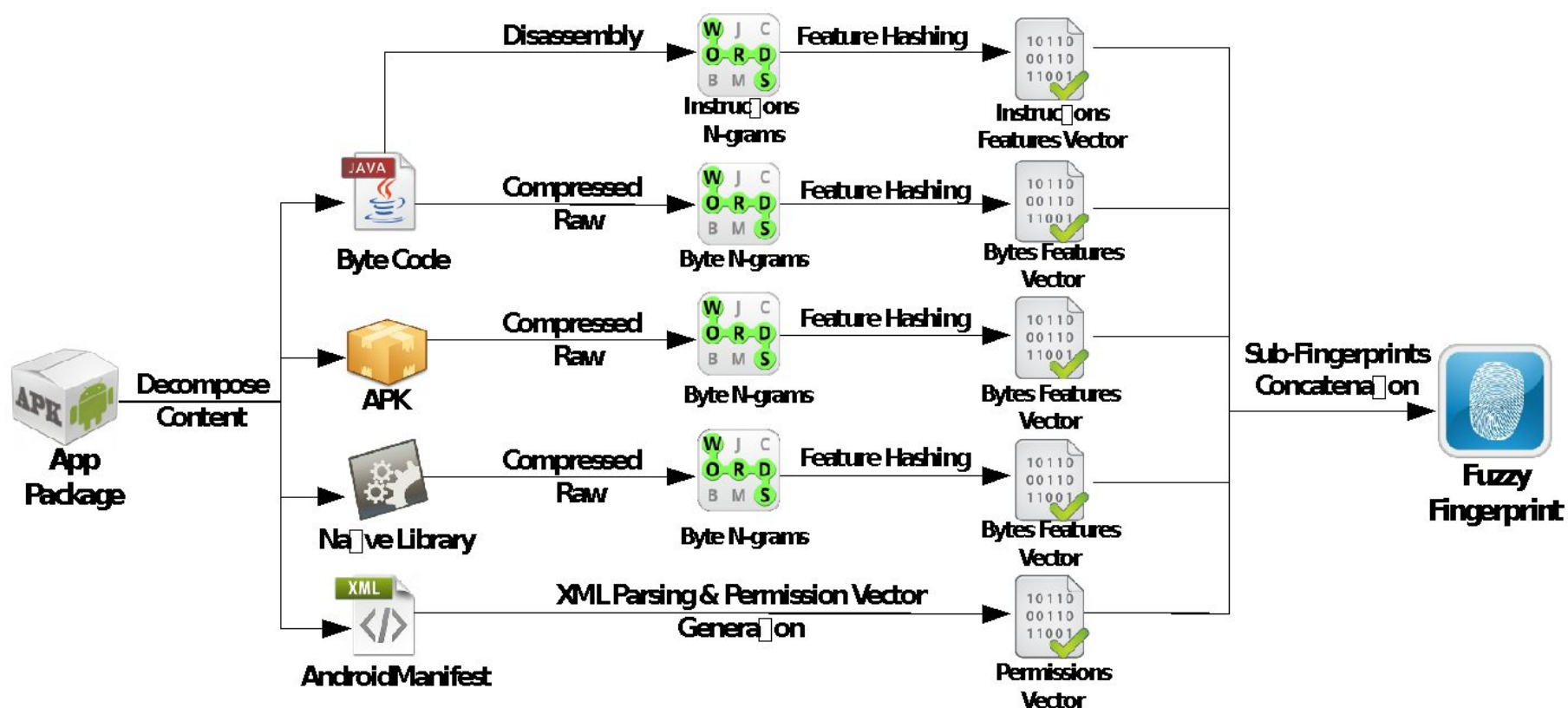
$$Jaccard\_bitwise(A, B) = \frac{Ones(A.B)}{Ones(A+B)}$$
$$0 \leq Jaccard\_bitwise(X, Y) \leq 1$$

A = 1011 0110 1001 0001  
B = 0011 1010 1001 0000  
A+B = 1011 1110 1001 0001  
A.B = 0011 0010 1001 0000

Ones(A+B) = 9  
Ones(A.B) = 5  
Jaccard(A,B) = 5 / 9

# Methodology Overview

16





# APK-DNA Fingerprint Similarity

17

- Compare APK-DNA Similarities.

---

## Algorithm 2: APK-DNA Similarity Computation

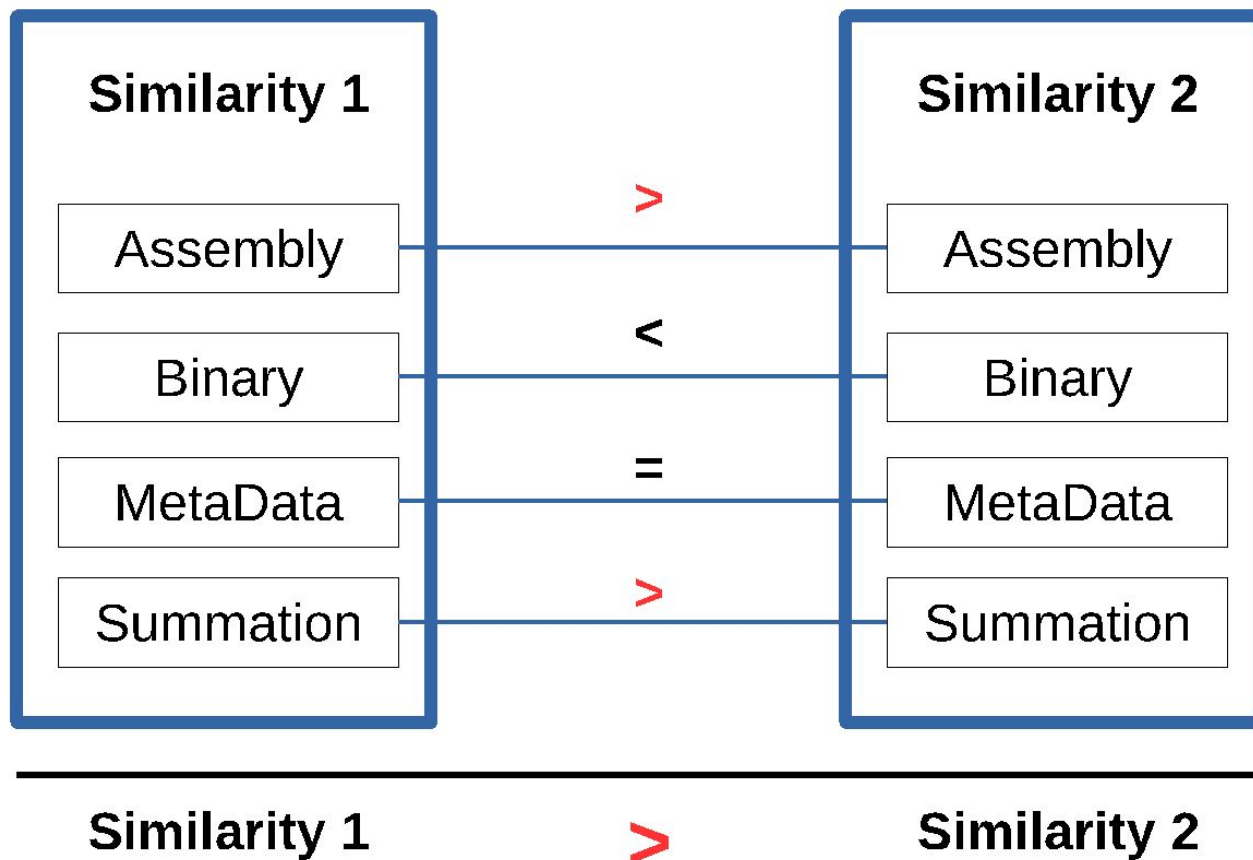
---

```
input  : APK-DNA A : list
        APK-DNA B : list
output: similarity-list: list
similarity-list = empty-list();
for content in content categories do
    | similarity = Jaccard_bitwise(A[content],B[content]) ;
    | similarity-list.add(similarity);
end
summation = sum(similarity-list);
similarity-list.add(summation);
```

---

# APK-DNA Fingerprint Similarity

18



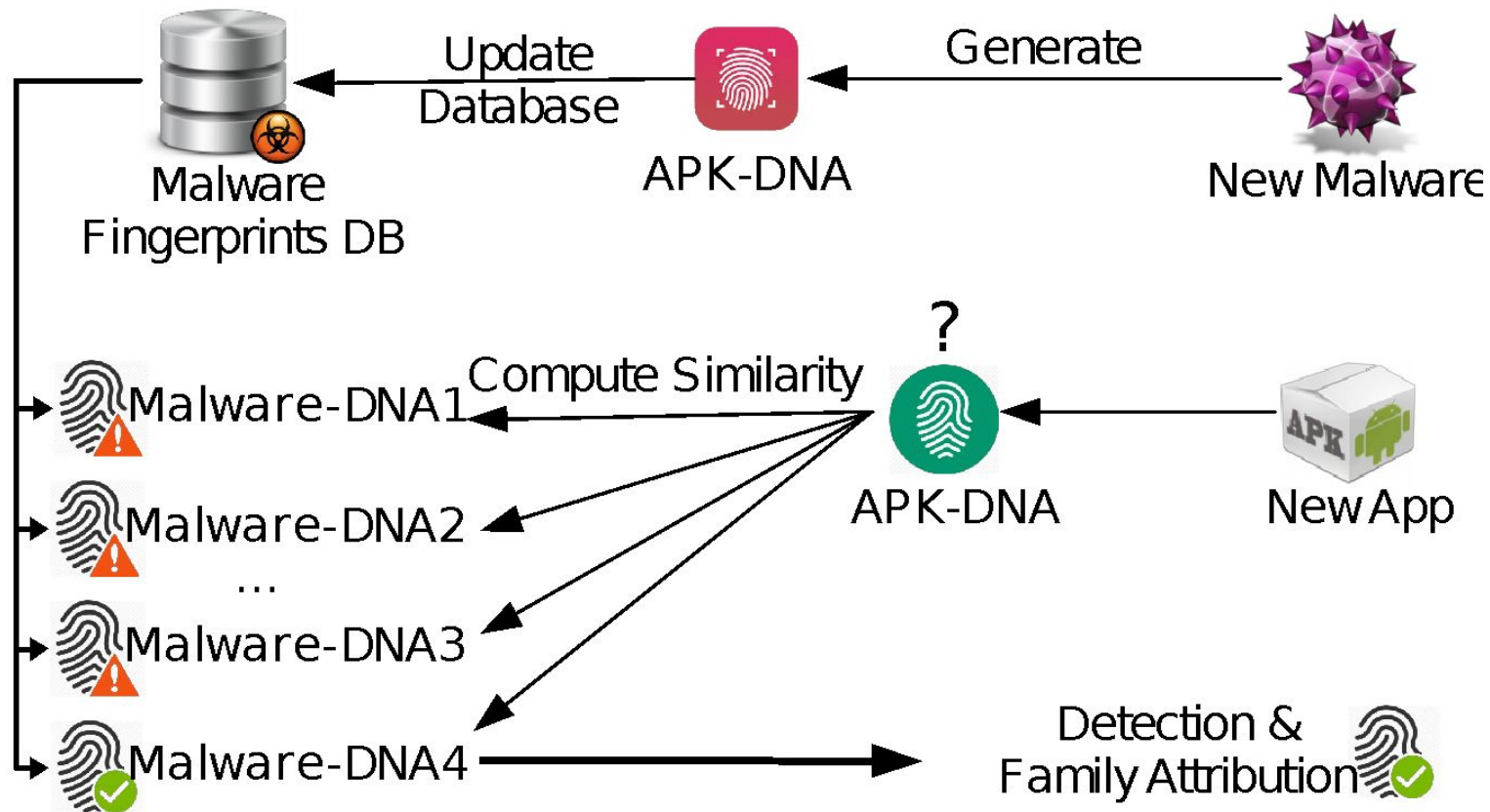
# Framework Approaches

19

- Peer-Matching
- Family-Fingerprint

# Peer-Matching approach

20



# Family-Fingerprinting Approach

21

- Compute
- Overview

# Family-Fingerprinting Compute

22

---

## Algorithm 4: Family Fingerprint Computation

---

input : M alware Family X Fingerprints: Set

output: Family X Fingerprint: FP\_X

FP\_X = new bitvector[Zeros];

for *fprint* in *Fingerprints* do

    FP\_X {meta} = FP\_X {meta} or fprint {meta};

    FP\_X {bin} = FP\_X {bin} or fprint {bin};

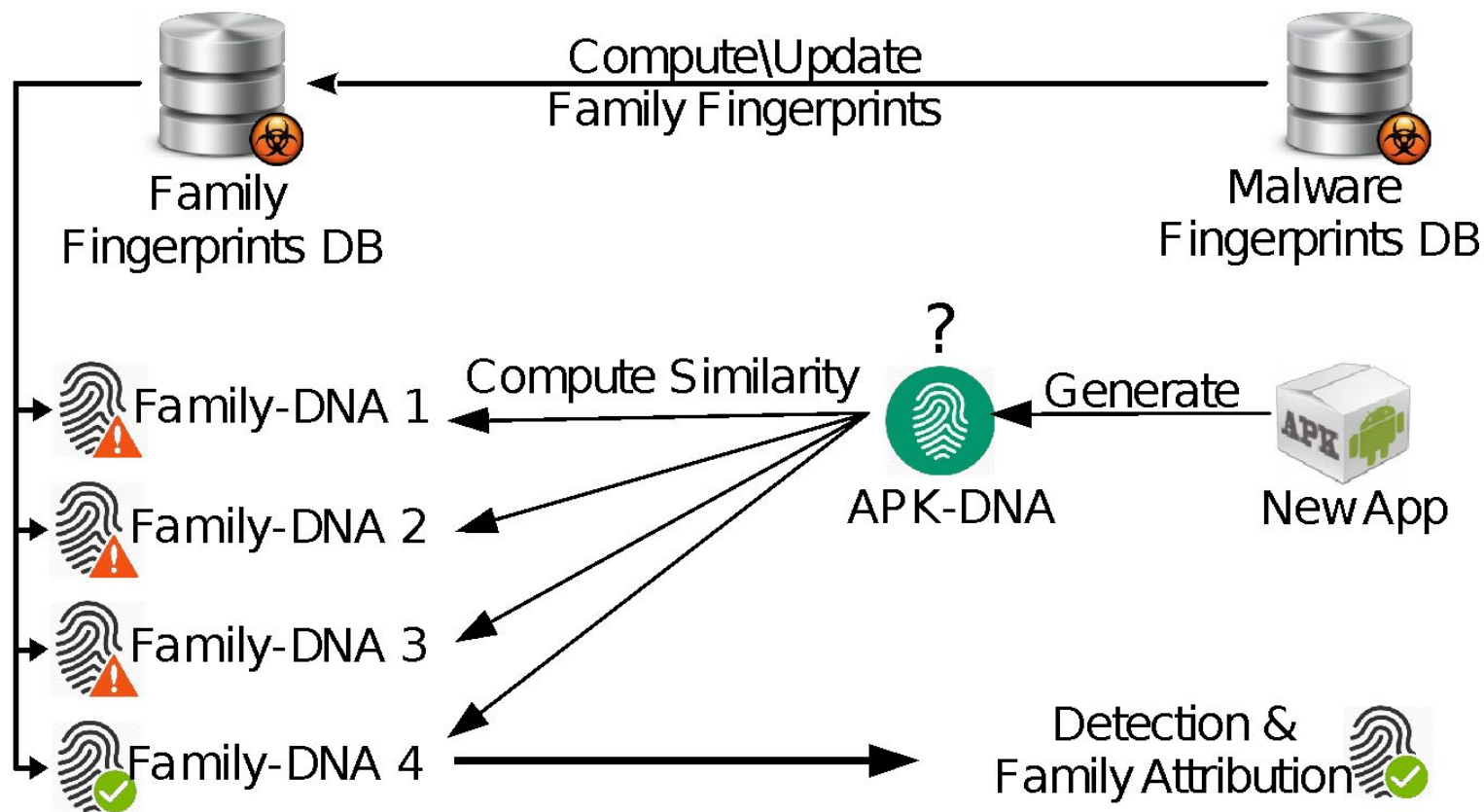
    FP\_X {asm} = FP\_X {asm} or fprint {asm};

end

---

# Family-Fingerprinting Overview

23



# Evaluation

24

- Android Malware Dataset.
- Evaluation Metrics.



# Evaluation Dataset

25

## □ Android Malware Genome Project Dataset

#	Malware Family/ Apps	Number of Samples
0	AnserverBot	187
1	KMin	52
2	DroidKungFu4	96
3	GoldDream	47
4	Geinimi	69
5	BaseBridge	122
6	DroidDreamLight	46
7	DroidKungFu3	309
8	Benign Apps	100

# Evaluation Metrics

26

$$Precision = \frac{TP}{TP+FP}$$

$$Recall = \frac{TP}{TP+FN}$$

$$F1 - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

# Results

27

- Family-Fingerprint Approach
- Peer-Matching Approach
- Voting Similarity vs Merging Similarity.

# Results (Family-Fingerprint)

28

Fingerprint Setup	<i>F1-Score</i>	<i>Precision</i>	<i>Recall</i>
Assembly	69%	88%	68%
APK	33%	36%	32%
Permission	69%	84%	70%
Dex	41%	46%	43%
Assembly, Permission, Dex, APK	81%	88%	80%
Assembly, Permission	82%	88%	81%
<b>Assembly, Permission, Dex</b>	<b>85%</b>	<b>89%</b>	<b>84%</b>
<b>Best Fingerprint Setup</b>	<b>85%</b>	<b>89%</b>	<b>84%</b>

Table 2: Accuracy Results of the Family-Fingerprinting Approach

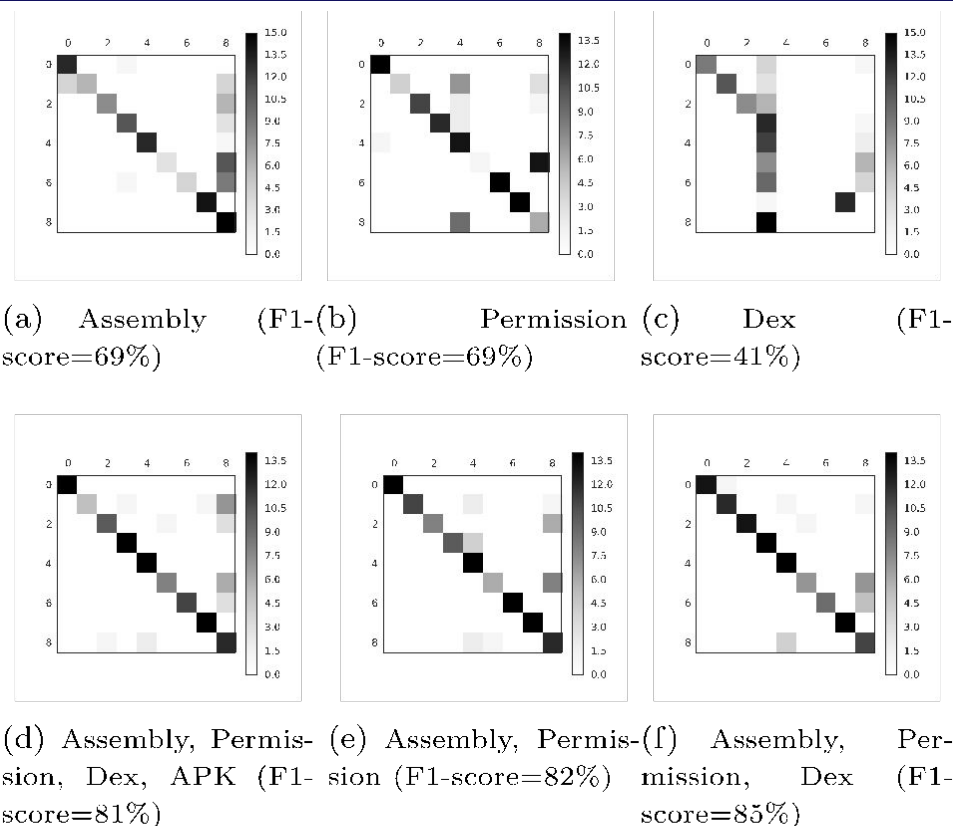


Figure 8: Confusion Matrices of Family-Fingerprint Approach for each Fingerprint Setup

# Results (Peer-Matching)

29

Fingerprint Setup	<i>F1-Score</i>	<i>Precision</i>	<i>Recall</i>
Assembly	91%	91%	90%
Apk	46%	48%	44%
Permission	81%	82%	80%
Dex	80%	90%	84%
Assembly, Permission, Dex, APK	84%	91%	81%
Assembly, Permission, Dex	93%	94%	93%
<b>Assembly, Permission</b>	<b>94%</b>	<b>95%</b>	<b>94%</b>
<b>Best Fingerprint Setup</b>	<b>94%</b>	<b>95%</b>	<b>94%</b>

Table 3: Accuracy Result of Peer-Matching Approach

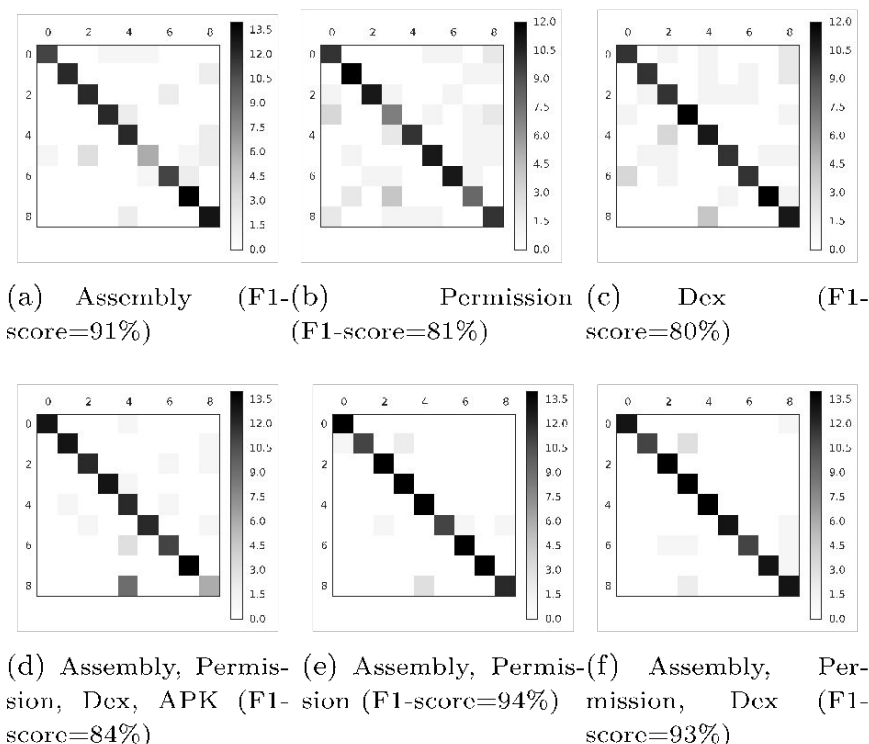


Figure 9: Confusion Matrices of Peer-Matching Approach for each Fingerprint Setup

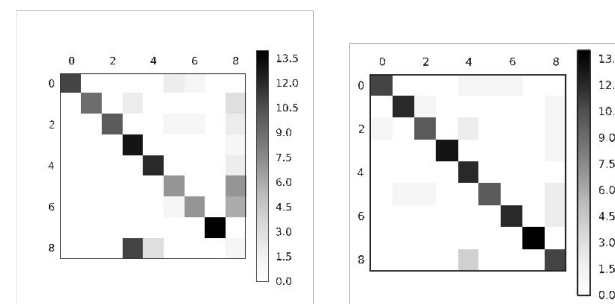
# Results (Voting VS Merging)

30

## □ Detection Performance using Peer-Voting Vs Merged Similarity.

Fingerprint Setup	<i>F1-Score</i>	<i>Precision</i>	<i>Recall</i>
Merged in Family-Approach	72%	84%	72%
<b>Peer-Voting in Family-Approach</b>	<b>85%</b>	<b>89%</b>	<b>84%</b>
Merged in Peer-Approach	86%	87%	86%
<b>Peer-Voting in Peer-Approach</b>	<b>94%</b>	<b>95%</b>	<b>94%</b>

Table 4: Accuracy Result of ROAR Using Merged Fingerprint



(a) Merged Fingerprint in Family-Fingerprint Approach (F1-score=72%) (b) Merged Fingerprint in Peer-Matching Approach (F1-score=86%)

Figure 10: Confusion Matrices Of ROAR Approaches Using Merged Fingerprint

# Summary

31

- ❑ Comprehensive fuzzy fingerprinting design for investigating Android malware variations, APK-DNA.
- ❑ Captures not only the binary of the APK file, but also both its structure and semantics.
- ❑ Framework for Android malware detection Using APK-DNA, peer-matching and family-fingerprint approaches.
- ❑ The evaluation demonstrated very promising results

# Acknowledgement

32

- ❑ Research members of the NCFTA Canada lab headed by Prof. Mourad Debbabi
- ❑ The anonymous reviewers and Timothy Vidas for their insightful comments.
- ❑ The DFRWS USA 2016 organization comity.



# Questions

33

- Thank you

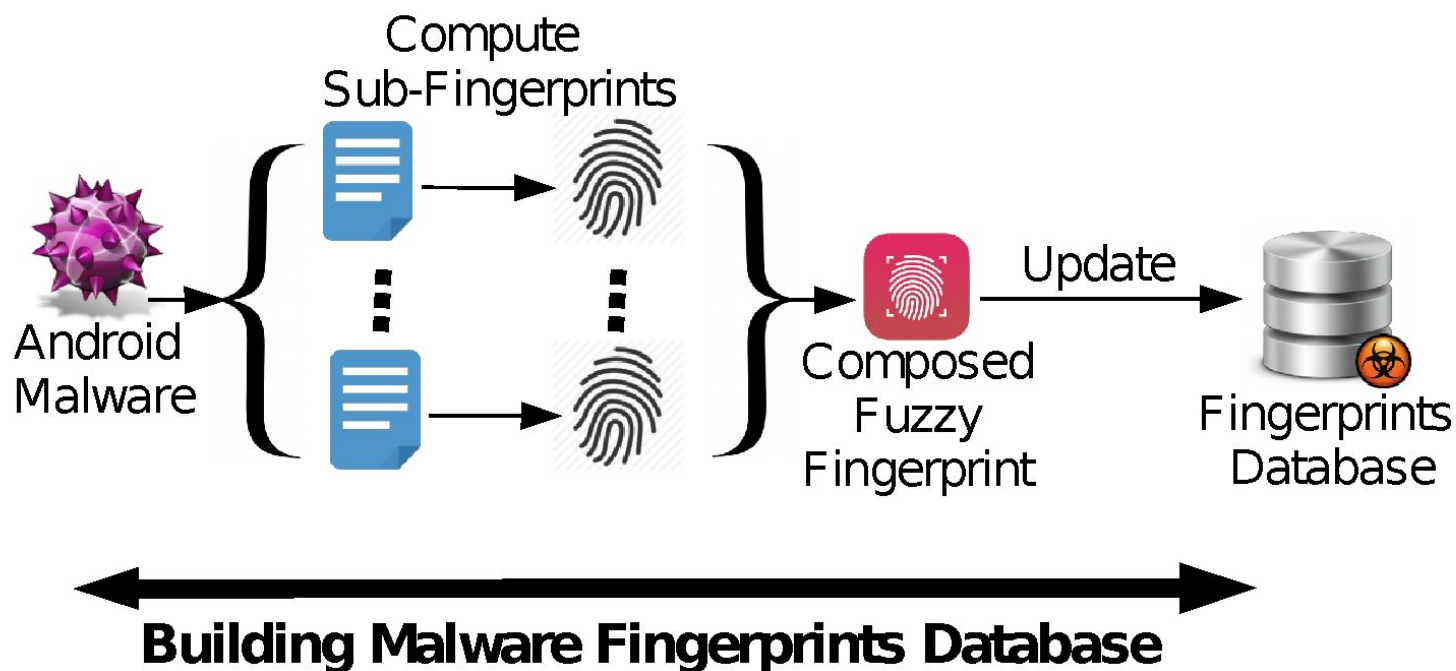
ElMouatez Billah Karbab  
Ph.D. Student  
Cyber Security R&D

[e\\_karbab@encs.concordia.ca](mailto:e_karbab@encs.concordia.ca)  
[mtz@ncfta.ca](mailto:mtz@ncfta.ca)

# Appendix

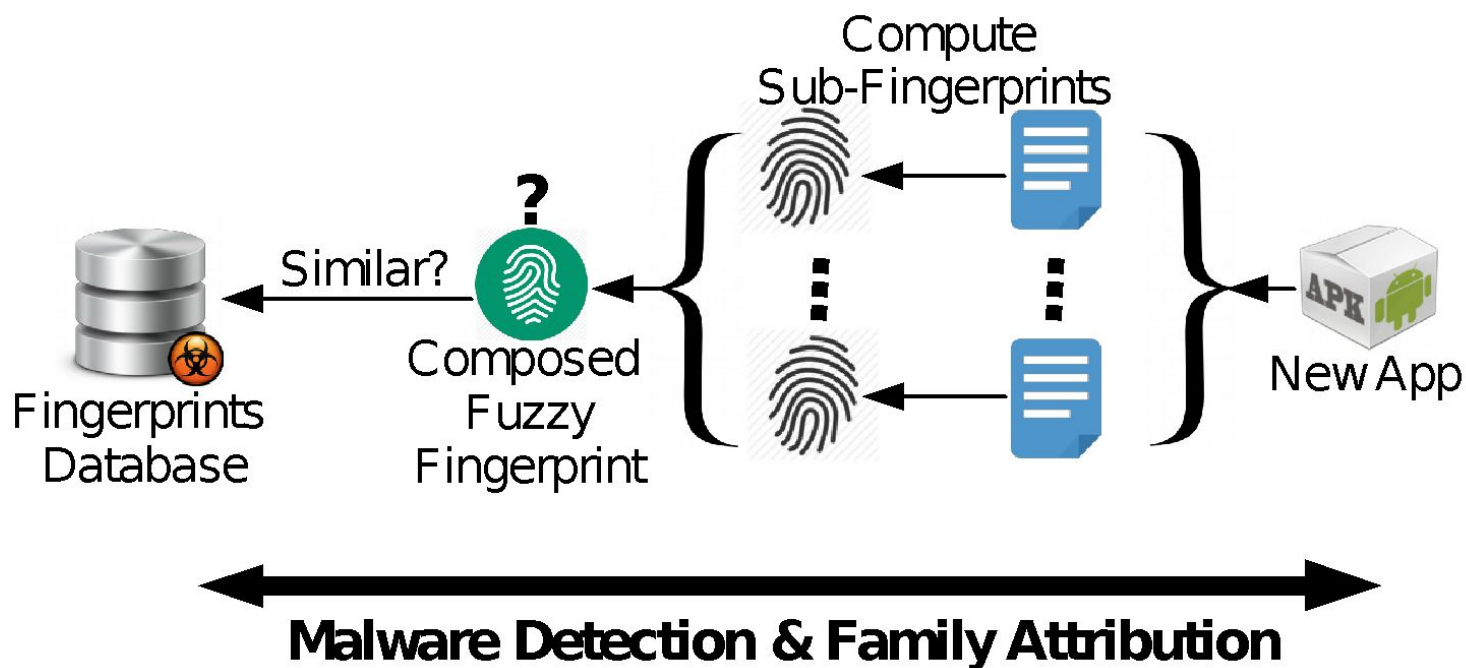
# Approach (cont'd)

35



# Approach (cont'd)

36



# APK-DNA Fingerprint Similarity

37

---

## Algorithm 3: Peer-Fingerprint Voting Mechanism

---

```
input  : similarity-list A-B: list
        similarity-list A-C: list
output: Decision
A-B-count = 0 ;
A-C-count = 0 ;
for content in content categories do
    | if  $A-B[content] > A-C[content]$  then
    |   | A-B-count += 1;
    | else
    |   | A-C-count += 1;
    | end
end
if  $A-B-count > A-C-count$  then
    | Decision = A-B;
else
    | Decision = A-C;
end
```

---