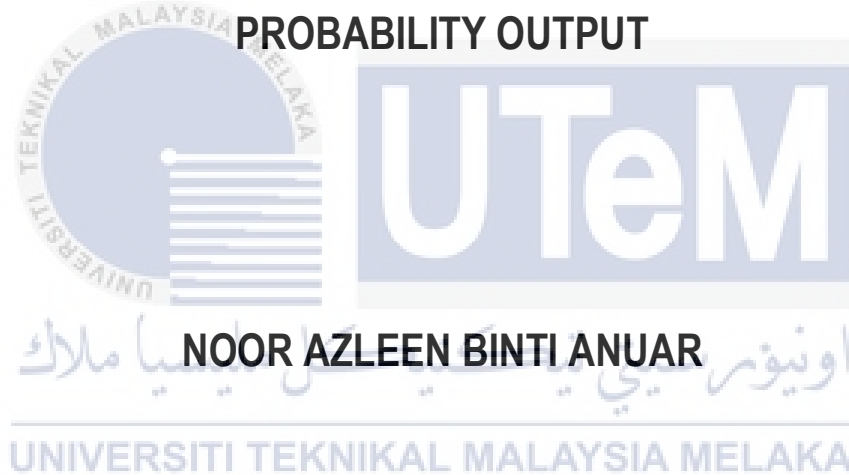




**A MOBILE MALWARE DETECTION FRAMEWORK BASED ON  
ENSEMBLE CLASSIFIER OF MULTIPLE N-GRAM OPCODE  
PROBABILITY OUTPUT**



**MASTER OF SCIENCE IN INFORMATION AND  
COMMUNICATION TECHNOLOGY**

**2023**



**Faculty of Information and Communication Technology**

A large, faded version of the UTeM logo is centered in the background of the page, behind the title text.

**A MOBILE MALWARE DETECTION FRAMEWORK BASED ON  
ENSEMBLE CLASSIFIER OF MULTIPLE N-GRAM OPCODE  
PROBABILITY OUTPUT**

اونيورسي تيكنيكل مليسيا ملاك  
UNIVERSITI TEKNIKAL MALAYSIA MELAKA

**NOOR AZLEEN BINTI ANUAR**

**Master of Science in Information and Communication Technology**

**2023**

**A MOBILE MALWARE DETECTION FRAMEWORK BASED ON ENSEMBLE  
CLASSIFIER OF MULTIPLE N-GRAM OPCODE PROBABILITY OUTPUT**

**NOOR AZLEEN BINTI ANUAR**



**Faculty of Information and Communication Technology**

**UNIVERSITI TEKNIKAL MALAYSIA MELAKA**

**2023**

## DECLARATION

I declare that this thesis entitled “A Mobile Malware Detection Framework based on Ensemble Classifier of Multiple N-Gram Opcode Probability Output” is the result of my own research work except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.

 Signature : 

Name : Noor Azleen Binti Anuar  
Date : 18th August 2023

اونيورسيٲى ٲكنيكل ماليسيا ملاك

---

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

## APPROVAL

I hereby declare that I have read this thesis and in my opinion this thesis is sufficient in terms of scope and quality for the award of Master of Science in Information and Communication Technology.

 Signature : .....  
Supervisor Name : Ts. Dr. Mohd Zaki Bin Mas'ud  
Date : 18th August 2023  
  
اوتيمر سیتی تکنیکل ملیسیا ملاک  
UNIVERSITI TEKNIKAL MALAYSIA MELAKA

## DEDICATION

This thesis is dedicated with  
Deepest love and affections to my beloved parents,

**Anuar bin Jabbar and Rudziah binti Adam**

Brother and sisters

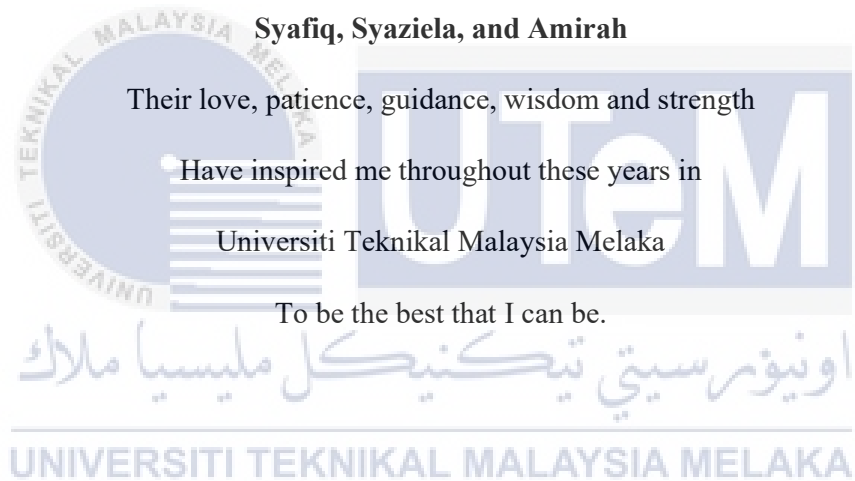
**Syafiq, Syaziela, and Amirah**

Their love, patience, guidance, wisdom and strength

Have inspired me throughout these years in

Universiti Teknikal Malaysia Melaka

To be the best that I can be.



## ABSTRACT

The advancement of mobile devices nowadays lets users do varieties of activities including surfing the internet, online banking transactions, engaging in social networking and hence increasing the usage of mobile devices. This scenario opens the possibility for cybercriminals to launch a mobile malware attack towards users. The complexity of detecting mobile malware also contributes to the possibility of mobile malware remaining dormant in the application store which can expose users to being tricked into installing the infected programs. Current mobile malware detection methods such as Static analysis and signature-based detection can address these issues, but it can be very difficult to detect zero-day or obfuscated code because it relies on a unique signature. Meanwhile, Dynamic analysis and anomaly-based detection can curb the problem, yet it can result in a relatively high rate of false alerts. In addition, a single model classifier is not strong enough to produce a good detection result. Based on this reason, this research intends to enhance the current Mobile Malware Detection Framework using multiple N-Gram opcode probability output and weighted ensemble to enhance the accuracy, TPR, and FPR. The aim of this research is to identify the features of malicious activity from mobile malware application through static analysis. The features obtained were used in formulating and evaluating the enhanced MMD Framework. The generation of N-Gram opcode sequence represents the malicious features and feature selection method is used to search for optimum features. Additionally, the weighted ensemble method is introduced to combine several probability outputs from multiple classification models. Particle Swarm Optimization is used in searching for optimum weight to be used together with the probability output to improve mobile malware detection. In conclusion, the proposed MMD Framework had shown an enhanced performance with an accuracy of 96.55%, TPR of 99.10%, and FPR of 0.90%. Based on the encouraging results, future studies could explore the possibility of using a dynamic analysis detection approach and applying n-gram to features other than opcode sequence. Ultimately, other datasets and other mobile malware variants should also be explored in future.

**RANGKA KERJA PENGESANAN PERISIAN HASAD MUDAH ALIH  
BERDASARKAN PENGELAS ENSEMBEL BAGI OUTPUT KEBARANGKALIAN  
OPKOD N-GRAM BERBILANG**

**ABSTRAK**

*Kemajuan peranti mudah alih pada masa kini membolehkan pengguna melakukan pelbagai aktiviti termasuk melayari internet, transaksi perbankan dalam talian, terlibat dalam rangkaian sosial dan lain-lain. Walau bagaimanapun, senario ini membuka kemungkinan penjenayah siber melancarkan serangan perisian hasad mudah alih terhadap pengguna. Kerumitan dalam mengesan perisian hasad mudah alih juga menyumbang kepada kemungkinan perisian hasad mudah alih kekal tidak aktif di kedai aplikasi yang boleh mendedahkan pengguna untuk ditipu untuk memasang program yang dijangkiti. Kaedah pengesanan perisian hasad mudah alih semasa seperti analisis statik dan pengesanan berasaskan tandatangan boleh menangani isu-isu ini, tetapi ia boleh menjadi sangat sukar untuk mengesan kod sifar hari atau obfuscated kerana ia bergantung pada tandatangan yang unik. Sementara itu, analisis dinamik dan pengesanan berasaskan anomali dapat membendung masalah, namun ia boleh mengakibatkan kadar amaran palsu yang agak tinggi. Di samping itu, pengelas model tunggal tidak cukup kuat untuk menghasilkan hasil pengesanan yang baik. Berdasarkan sebab ini, penyelidikan ini berhasrat untuk meningkatkan kaedah Pengesanan Perisian Hasad Mudah Alih semasa menggunakan pelbagai ciri urutan opkod N-Gram dan kaedah klasifikasi ensemble berwajaran dari segi ketepatan pengesananannya, Kadar Positif Sebenar dan Kadar Positif Palsu. Matlamat penyelidikan ini adalah untuk mengenal pasti ciri aktiviti berniat jahat daripada aplikasi perisian hasad mudah alih melalui analisis statik. Ciri-ciri yang diperolehi digunakan dalam merumus dan menilai Rangka Kerja MMD yang dipertingkatkan. Penjanaan jujukan opkod N-Gram mewakili ciri berniat jahat dan kaedah pemilihan ciri digunakan untuk mencari ciri optimum. Selain itu, kaedah ensemble berwajaran diperkenalkan untuk menggabungkan beberapa keluaran kebarangkalian daripada pelbagai model klasifikasi. Pengoptimuman Kawanan Zarah digunakan dalam mencari berat optimum untuk digunakan bersama-sama dengan output kebarangkalian untuk meningkatkan pengesanan perisian hasad mudah alih. Kesimpulannya, Rangka Kerja MMD yang dicadangkan telah menunjukkan prestasi yang dipertingkatkan dengan ketepatan sebanyak 96.55%, TPR 99.10%, dan FPR 0.90%. Berdasarkan keputusan yang menggalakkan, kajian masa depan boleh meneroka kemungkinan menggunakan pendekatan pengesanan analisis dinamik dan menggunakan n-gram pada ciri selain daripada urutan opcode. Akhirnya, set data lain dan varian perisian hasad mudah alih yang lain juga harus diterokai pada masa hadapan.*



## ACKNOWLEDGMENT

In the name of Allah, the Most Gracious and the Most Merciful

First and foremost, All Praise to Allah SWT the Almighty for the strength and blessing in completing my thesis and Masters journey. A sincere gratitude I'll give to my supervisor, Ts. Dr. Mohd Zaki bin Mas'ud, whose expertise, understanding, patience, and guidance had helped me enhance my graduate experience. I am so grateful for his priceless effort in helping me whenever I find difficulties in completing my task and giving comments for improvement of this thesis.

Special thanks to my co-supervisor; Ts. Dr. Nazrulazhar bin Bahaman, my mentor; Ts. Nor Azman bin Mat Ariff for their support and aid in making my Masters journey a success. My appreciation also goes to Universiti Teknikal Malaysia Melaka (UTeM) for the opportunity given.

Last but not least, from the bottom of my heart a highest gratitude to my family for their love and support. Especially to my father, Anuar bin Jabbar and my mother, Rudziah binti Adam for their encouragements and blessings. Finally, to those who indirectly contributed to this research, your kindness has inspired me to embark on this journey.

## TABLE OF CONTENTS

	PAGE
<b>DECLARATION</b>	
<b>APPROVAL</b>	
<b>DEDICATION</b>	
<b>ABSTRACT</b>	<b>i</b>
<b>ABSTRAK</b>	<b>ii</b>
<b>ACKNOWLEDGEMENTS</b>	<b>iii</b>
<b>TABLE OF CONTENTS</b>	<b>iv</b>
<b>LIST OF TABLES</b>	<b>vii</b>
<b>LIST OF FIGURES</b>	<b>viii</b>
<b>LIST OF APPENDICES</b>	<b>x</b>
<b>LIST OF ABBREVIATIONS</b>	<b>xi</b>
<b>LIST OF SYMBOLS</b>	<b>xiv</b>
<b>LIST OF PUBLICATIONS</b>	<b>xv</b>
<b>CHAPTER</b>	
<b>1. INTRODUCTION</b>	<b>1</b>
1.1 Background	1
1.2 Research Problems	3
1.3 Research Questions	6
1.4 Research Aim and Objectives	8
1.5 Research Scope	10
1.6 Research Contributions	11
1.7 Thesis Organization	12
1.8 Summary	15
<b>2. LITERATURE REVIEW</b>	<b>16</b>
2.1 Introduction	16
2.2 Chapter Objective	16
2.3 Chapter Outline	17
2.4 Overview of Mobile Malware Issues	18
2.4.1 Mobile Malware Evolution	18
2.4.2 Mobile Malware Threats and Impact	20
2.4.3 Android Architecture and Security Framework	21
2.5 Mobile Malware Detection (MMD)	23
2.5.1 Mobile Malware Analysis Approach	25
2.5.2 Mobile Malware Audit Data Source	27
2.5.3 Data Acquisition	30
2.5.4 Mobile Malware Detection Technique	32
2.5.5 Feature Selection Process	34
2.5.6 Mobile Malware Detection Classification Analysis	37
2.6 Classifier Algorithm	40
2.6.1 Support Vector Machine (SVM)	42
2.6.2 Single Classifier Technique	45

2.6.3	Multiple Classifier Technique	45
2.6.4	Taxonomy of Multiple Classifier Technique	47
2.6.5	Types of Multiple Classifier Technique	49
2.6.6	Combination Method	50
2.6.7	Weighted Ensemble Method using Particle Swarm Optimization (PSO)	53
2.7	Evaluation Process	55
2.8	Summary	56
<b>3.</b>	<b>RESEARCH METHODOLOGY</b>	<b>58</b>
3.1	Introduction	58
3.2	Chapter Objective	58
3.3	Chapter Outline	55
3.4	Research Process	60
3.4.1	Phase I - Information Gathering	61
3.4.2	Phase II – Analysing	62
3.4.3	Phase III – Synthesizing	62
3.4.4	Phase IV - Development and Testing	63
3.4.5	Phase V - Result Analysis	63
3.5	The Proposed MMD Framework	64
3.5.1	Phase I – Data Extraction	65
3.5.2	Phase II – Feature Vector Generation	68
3.5.3	Phase III – Machine Learning Classification and Ensemble Method	70
3.5.4	Phase IV – Evaluation of Proposed Framework	71
3.6	Mapping of Research Methodology and Research Objectives	72
3.7	Summary	72
<b>4.</b>	<b>MOBILE MALWARE BEHAVIOUR THROUGH N-GRAM OPCODE SEQUENCE</b>	<b>74</b>
4.1	Introduction	74
4.2	Chapter Objective	74
4.3	Chapter Outline	75
4.4	Mobile Malware Behaviour Experimental Approach Overview	76
4.5	Mobile Malware Behaviour Analysis	77
4.5.1	Remote Access Trojan Behaviour Analysis	79
4.5.2	Banking Trojan Behaviour Analysis	80
4.5.3	Adware Behaviour Analysis	81
4.5.4	Discussion on Malicious Mobile Malware Behaviour	82
4.6	Mobile Malware Behaviour through Opcode Sequence	84
4.6.1	Capturing Device Information	85
4.6.2	Connection with External Server	86
4.6.3	Execution of Malicious Payload	87
4.7	N-Gram Opcode Sequence Generation	88
4.8	Summary	91
<b>5.</b>	<b>WEIGHTED ENSEMBLE METHOD USING MULTIPLE N-GRAM OPCODE PROBABILITY OUTPUT</b>	<b>92</b>
5.1	Introduction	92

5.2	Chapter Objective	93
5.3	Chapter Outline	93
5.4	Feature Selection Evaluation Process	94
5.4.1	Feature Selection Method Evaluation Experiment	94
5.4.2	Feature Selection Method Evaluation and Analysis Results	95
5.5	The Proposed N-gram Opcode Sequence and Weighted Ensemble Method in MMD	96
5.5.1	Result for BOW using PSO	97
5.5.2	Result for CS using PSO	99
5.5.3	Result for IG using PSO	101
5.6	Ensemble of Multiple N-Gram Opcode Probability Output using PSO in Mobile Malware Detection	103
5.7	Summary	106
<b>6.</b>	<b>CONCLUSION AND RECOMMENDATIONS</b>	<b>107</b>
6.1	Introduction	107
6.2	Chapter Objective	108
6.3	Research Contribution	111
6.3.1	Mobile Malware Behaviour Analysis	112
6.3.2	N-Gram Opcode Sequence	114
6.3.3	Weighted Ensemble Method	110
6.3.4	Improved Effectiveness of Mobile Malware Detection	115
6.3.5	Summary of Research Contributions	117
6.4	Recommendation and Future Works	118
6.5	Conclusion	119
	<b>REFERENCES</b>	<b>120</b>
	<b>APPENDICES</b>	<b>135</b>

## LIST OF TABLES

TABLE	TITLE	PAGE
1.1	Research problems	5
1.2	Summary of research questions	7
1.3	Summary of research problems (RP), research questions (RQ) and research objectives (RO)	9
1.4	Summary of research contributions	11
2.1	Implementation of analysis approach	26
2.2	Audit data source overview	29
2.3	Techniques used in feature selection	35
2.4	The advantages and disadvantages of each MMD element	38
2.5	Application of multiple classifier techniques by previous researchers	49
3.1	Research methodology mapping with RO 1, RO 2, and RO 3	72
4.1	Mobile malware types and its malicious activities	76
4.2	List of the encoded N-gram opcode sequence	89
4.3	Sample of the opcode sequence encoding scheme	90
4.4	List of the encoded n-gram opcode sequence	90
5.1	Classification accuracy results	95
5.2	Result for weighted ensemble method using PSO for BOW	97
5.3	Result for weighted ensemble method using PSO for CS	99
5.4	Result for weighted ensemble method using PSO for IG	101
5.5	Performance evaluation for proposed method	105
6.1	Comparison for the results obtained from the proposed MMD framework with previous researchers	115

## LIST OF FIGURES

FIGURE	TITLE	PAGE
1.1	Mobile malware attack trends (Kaspersky, 2022)	2
1.2	Evasion techniques (McAfee Labs, 2017)	3
1.3	Thesis outline	17
2.1	The overview of chapter two	16
2.2	Timeline of mobile malware evolution	20
2.3	Architecture of Android platform (Yerima et al., 2013)	22
2.4	Android malware detection taxonomy	34
2.5	Observations for supervised and unsupervised learning algorithm	41
2.6	SVM hyperplane	43
2.7	Basic scheme for PSO algorithm (Marini and Walczak, 2015)	55
3.1	Outline for chapter three	59
3.2	Overall research process	60
3.3	The proposed enhanced mobile malware detection framework	64
3.4	Steps involved in mobile malware features through opcode sequence	67
3.5	Steps involved in N-Gram generator	68
3.6	Steps involved in classification algorithm and weighted ensemble using PSO	71
4.1	Overview of chapter four	75
4.2	Opcode Traces for Remote Access Trojan	79
4.3	Opcode Traces for Banking Trojan	81
4.4	Opcode Traces for Adware	82
4.5	Malicious mobile malware behaviour	82
4.6	Malicious intention found during analysis	83
4.7	Opcode sequence used to retrieve information in a file	85
4.8	Opcode sequence used in information gathering	86
4.9	Opcode sequence used to connect to an external server	86

4.10	Opcode sequence used to gain root access	87
4.11	Opcode sequence used to access system file	87
4.12	Feature vector generation process	88
5.1	Chapter five outline	94
5.2	Classification accuracy using single classifier	96
5.3	BOW accuracy for each particles	99
5.4	CS accuracy for each particles	101
5.5	IG accuracy for each particles	103
5.6	Example for the implementation of machine learning classification and ensemble method phase	104
5.7	Comparison of accuracy for each weighted ensemble method	107
6.1	The objectives and contributions mapping	112
6.2	Mapping of research objectives, activities and outcomes	117



## LIST OF APPENDICES

APPENDIX	TITLE	PAGE
A	Implementation of MCT by previous researchers	135
B	Flowchart of the whole system	136
C	Opcode representation	137





## LIST OF ABBREVIATIONS

AB	-	Anomaly Based
AMD	-	Android Malware Dataset
API	-	Application Program Interface
APK	-	Android Package
AUC	-	Area Under Curve
BOW	-	Bag of Words
C&C	-	Command and Control
CPU	-	Central Processing Unit
CS	-	Chi-Square
DVM	-	Dalvik Virtual Machine
EA	-	Evaluation Algorithm
FN	-	False Negative
FP	-	False Positive
FPR	-	False Positive Rate
GA	-	Genetic Algorithm
GB	-	Gigabyte
GPS	-	Global Positioning System
GSS	-	Greedy Search Strategy
HIS	-	Hybrid Intelligent System
HTTP	-	Hypertext Transfer Protocol
ICCID	-	Integrated Circuit Card ID

IDS	-	Intrusion Detection System
IG	-	Information Gain
IMEI	-	International Mobile Station Equipment Identity
IMSI	-	International Mobile Subscriber Identity
IP	-	Internet Protocol
IPC	-	Inter Process Communication
kNN	-	K-Nearest Neighbour
MCT	-	Multiple Classifier Technique
MMD	-	Mobile Malware Detection
OS	-	Operating System
PO	-	Probability Output
POC	-	Proof of Concept
PSO	-	Particle Swam Optimization
RAM	-	Random Access Memory
RC	-	Research Contribution
RO	-	Research Objective
RP	-	Research Problem
RQ	-	Research Question
SB	-	Signature Base
SBS	-	Sequential Backward Search
SEO	-	Search Engine Optimization
SFS	-	Sequential Forward Search
SMS	-	Short Message Services
SU	-	Symmetrical Uncertainty
SVM	-	Support Vector Machine
TN	-	True Negative
TNR	-	True Negative Rate

TP	-	True Positive
TPR	-	True Positive Rate
UID	-	Unique Identifier



## LIST OF SYMBOLS

$f_{ngram}$	-	Normalize occurrence of n-gram opcode sequence frequency
$f_c$	-	Normal n-gram frequency
$min(f)$	-	The minimum value of the attribute
$max(f)$	-	The maximum value of attribute
$S$	-	The scaling factor for the output range
$T$	-	The translation factor for the output range
$w_i$	-	Weight Assigned
$h_i$	-	Probability Output for each N-Gram Classifier
$h(x)$	-	Weighted Combination of Probability Output of N-Gram Classifiers

## LIST OF PUBLICATIONS

1. Noor Azleen Anuar, Mohd Zaki Mas'ud, Nazrulazhar Bahaman, and Nor Azman Mat Ariff., 2020. Mobile Malware Behavior through Opcode Analysis. *International Journal of Communication Networks and Information Security*, 12(3), pp. 345–354.
2. Noor Azleen Anuar, Mohd Zaki Mas'ud, Nazrulazhar Bahaman, and Nor Azman Mat Ariff., 2020. Analysis of Machine Learning Classifier in Android Malware Detection Through Opcode. *2020 IEEE Conference on Application, Information and Network Security (AINS)*, pp. 7–11. <https://doi.org/10.1109/AINS50155.2020.9315060>



## CHAPTER 1

### INTRODUCTION

#### 1.1 Background

With the great evolution of technology and communication, all living souls are obligated to use their mobile devices on a regular basis. Mobile gadgets enable people to communicate with one another from anywhere in the globe. Nonetheless, the presence of mobile devices provided enough opportunities for attackers to exploit them by introducing malware into the devices without the users' awareness (Varna and Visalakshi, 2020). Due to the difficulty in detecting mobile malware, they now have the capacity to go unnoticed in the app store (Tenenboim-Chekina et al., 2013), causing users to be deceived into installing the infected applications. Once the software is installed on the victims' devices, malicious malware penetration and replication begin, which can create major difficulties if left undetected.

Malware, often known as malicious software, is well-known for its tendency to disrupt computer software and hardware. Attackers frequently exploit it to take advantage of accessible resources and for other cybercriminal purposes such as stealing users' data, passwords, and credit card numbers. This assault might occur when victims open an infected email or download malicious software. When a malicious link or attachment received in an email is clicked, malware may be installed without the user's knowledge, and a ransomware attack may happen, causing the entire system to freeze. This might

potentially expose sensitive information, financial and corporate information. In addition, the attacker could obtain and extract the victim's login credentials and account information.

According to the I Threat Evolution in Q1 2020 issued by Kaspersky, (Kaspersky, 2022) the discovery of latest mobile malware variations had decreased in the year 2021. In Quarter 1, Quarter 2, and Quarter 3 of 2020, the number of attacks were 14,446,496, followed by 14,203,865, and 16,440,099. The highest number of attacks were in the Quarter 4 of 2020 with 18,085,657 of attacks. Following that, it can be seen that there was a decrease in the year 2021 with 15,239,031 number of attacks in Quarter 1 of 2021. Next, in Quarter 2, Quarter 3, and Quarter 4 of 2021 were 14,465,670, followed by 9,599,322, and 6,931,266 number of attacks. Although it can be seen that there was a slight decrease in the mobile malware threats, the number of attacks were still high. Figure 1.1 shows the trends of malware attacks on mobile devices.

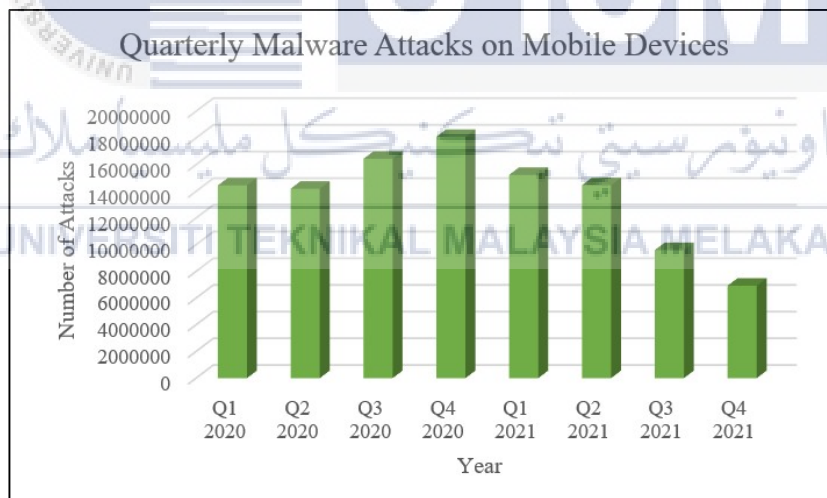


Figure 1.1: Mobile Malware Attack Trends (Kaspersky, 2022)

As reported by McAfee Labs (McAfee, 2017), several evasion techniques were used by mobile malware in their journey to detour into Android protection mechanisms. The most common evasion techniques are anti-sandbox, anti-security tools, code injection, anti-monitoring, and anti-debugging. Anti-sandbox can be defined as a method used to

isolate programs that is being executed in order to steer away from any unwanted applications to the system. Next, anti-security tools refer to a method used to bypass any detection done using security devices, or programs. On the other hand, code injection exists in the form of additional unwanted code or instructions to the binary in order to distract the disassembly view or waste the analyst's time. Anti-monitoring acts as a monitoring agent in order to prevent any reverse engineering from happening. Finally, yet importantly, anti-debugging is a program used so that the reverse engineering processing time will increase. All of the evasion techniques mentioned had made it harder to identify mobile malware. Figure 1.2 below shows the percentage of each evasion technique that had been reported in the year 2017.

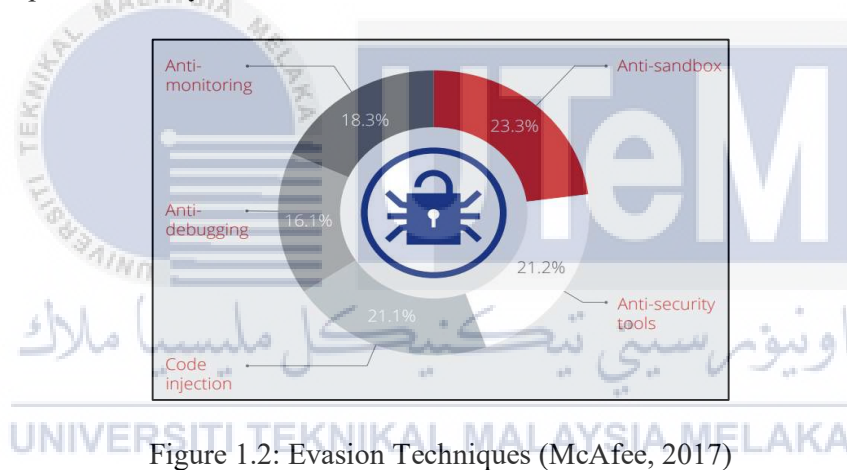


Figure 1.2: Evasion Techniques (McAfee, 2017)

## 1.2 Research Problems

Various analysis technique were implemented in order to solve the rising mobile malware threats. The two type of malware analysis are the static and dynamic analysis. Other than that, there are two types of classifier method in mobile malware detection that is single classifier and multiple classifier.

In static analysis, the code of an application is extracted and analyzed without having to be executed (Medvet and Mercaldo, 2016).The features of each application are