



**An Efficient Network Traffic Classification Based on
Vital Random Forest for High Dimensional Dataset**

by

ALHAMZA MUNTHER WARDI ALALOUSH

(1340211014)

A thesis submitted in fulfilment of the requirements for the degree
of Doctor of Philosophy

**School of Computer and Communication Engineering
UNIVERSITI MALAYSIA PERLIS**

2017

UNIVERSITI MALAYSIA PERLIS

DECLARATION OF THESIS

Author's full name : Alhamza Munther Wardi Alalousi

Date of birth : 03/12/1980.....

Title : An Efficient Network Traffic Classification Based on Vital.....
Random Forest for High Dimensional Dataset

.....

Academic Session : 2016/2017.....

I hereby declare that this thesis becomes the property of University Malaysia Perlis (UniMAP) and to be placed at the library of UniMAP. This thesis is classified as:

- CONFIDENTIAL** (Contains confidential information under the Official Secret Act 1972)
- RESTRICTED** (Contains where research was done) restricted Information as specified by the organization
- OPEN ACCESS** I agree that MY thesis is to be made immediately available as hard copy or online open access (full text)

I, the author, give permission to the UniMAP to reproduce this thesis in whole or in part for the purpose of research or academic exchange only (except during a period of — years, if so requested above).

Certified by:

SIGNATURE

SIGNATURE OF SUPERVISOR

(g1340811014/ A10833650)

(NEW IC NO. / PASSPORT NO.)

NAME OF SUPERVISOR

Date: _____

Date: _____

NOTES : * If the thesis is CONFIDENTIAL or RESTRICTED, please attach with the letter from the organization with period and reasons for confidentiality or restriction.

ACKNOWLEDGEMENTS

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
نَرْفَعُ دَرَجَاتٍ مِّنْ نَّشَأٍ وَفَوْقَ كُلِّ ذِي عِلْمٍ عَلِيمٌ ۖ ٧٦ (سورة يوسف)
((الْحَمْدُ لِلَّهِ الَّذِي بِنِعْمَتِهِ تَتِمُّ الصَّالِحَاتُ))

The favor, above all, before all, and after all, is entirely Allah's (SWT), to whom my never-ending thanks and praise are humbly due.

I would like to take this opportunity to convey my sincere thanks and deepest gratitude to my supervisor, Dr. Rozime Razif Bin Othman for all the help and valuable guidance provided to me during throughout my period of research and especially for his confidence in me. Also, for his insight and the precious scientific guidance, this greatly helped me in the research's progress and to accomplish of this thesis through his academic advice. I would also like to thank my co-supervisor Dr. Mohammed Anbar for his support, and invaluable suggestions during my research.

I would like to express my appreciation to Prof. Ir. Dr. Badlishah Ahmed the Dean of Computer and Communication Engineering School for his support and help towards my postgraduate affairs. My deep appreciation and special thanks go to the University Malaysia Perlis (UniMAP) and to the staff of the UniMAP for their cooperation, support and good treatment for foreign students; especially school members my colleagues. My Acknowledgement also goes to the Center of Graduate Studies, and the university library for their help and support.

Also, my deep appreciation and special thanks go to Prof. Dr. Syed Alwee Aljunid Bin Syed Junid the Dean of Research, Management and Innovation Centre for his unlimited helps and support during my study.

Moreover, I would like to thank those who are always in my heart; my father for his endless and continuous encouragement and constant support, my mother for her continuous prayers and inspiration. My sincere gratitude goes to my dearest brothers (Numan, Alabass) for their continuous supporting and encouragement. In addition to, my sisters for always keeping a smile on my face and motivating me all the time. My special acknowledgement goes to my cousin (Ammar) for his support and encouragement through my study journey. My thankful and grateful goes to Dr. Salah Noori for his sincere brotherly advices and continuous support during my study.

DEDICATION

((وَقُلْ رَبِّ ارْحَمْهُمَا كَمَا رَبَّيَانِي صَغِيرًا ۚ ۲۴))

سورة القصص

I would like to dedicate this thesis to the two persons (My father and my mother) who you always feel their presence beside me despite the long distance by introducing the continuous Doa'a, encourages and unlimited support during my study.

((قَالَ سَتَشِدُّ عَضُدَكَ بِأَخِيكَ وَنَجْعُلُ لَكُمَا سُلْطٰنًا ۚ ۳۵))

سورة الإسراء

I would also like to dedicate this thesis to my dearest brother Alabass for his exceptional sacrificing, continues help, prayers and support during my study.

((وصلى اللهم على سيدنا محمد وعلى اله وصحبه وسلم))

TABLE OF CONTENTS

	PAGE
DECLARATION OF THESIS	i
ACKNOWLEDGEMENTS	ii
TABLE OF CONTENTS	iv
LIST OF TABLES	x
LIST OF FIGURES	xiii
LIST OF ABBREVIATIONS	xvii
ABSTRAK	xx
ABSTRACT	xxi
CHAPTER 1 INTRODUCTION	
1.1 Introduction	1
1.2 Background	2
1.2.1 Network Traffic Engineering	2
1.2.2 Machine Learning	3
1.2.3 Original Random Forest (RF)	4
1.3 Problem Statement	5
1.4 Research Objective	8
1.5 Thesis Contributions	9
1.6 Research Methodology	10
1.7 Thesis Organization	12

CHAPTER 2 LITERATURE REVIEW

2.1	Introduction	13
2.2	Background	13
2.2.1	Network Traffic Classification	15
2.2.2	Features Selection Techniques	15
2.3	Related Works	17
2.3.1	Port-Based Network Traffic Classification	17
2.3.2	Signature Based Method	19
2.3.3	Behaviour-Based Classification	23
2.3.4	Network traffic classification using Machine learning	27
2.3.4.1	Naive Bayes	27
2.3.4.2	K-means Clustering	29
2.3.4.3	K-Nearest Neighbor	31
2.3.4.4	Probabilistic neural network	33
2.3.4.5	Support Vector Machine	36
2.3.4.6	Multi-class SVMs	40
2.3.4.7	Expectation Maximization	43
2.3.4.8	Hidden Markov Model (HMM)	45
2.3.4.9	C4.5 decision tree	50
2.3.4.10	Ensemble Classifier	52
2.4	Chapter Summary	58

CHAPTER 3 METHODOLOGY AND DESIGN

3.1	Introduction	59
3.2	Overview of Vital Random Forest (VRF)	59
3.3	Preprocessing Data in Data Mining	60
3.3.1	Feature Selection Techniques for Network Traffic	61
3.4	The Proposed VRF Phases	62
3.4.1	Phase 1: FSCGRV Based on Ranking Feature and Voting	62
3.4.1.1	Significance Feature Evaluator (SFE)	64
3.4.1.2	Fast Correlation-Based Filter	66
3.4.1.3	Chi-Squared feature selection	69
3.4.1.4	Gain Ratio Feature Selection	71
3.4.1.5	Fast Feature Ranking	73
3.4.1.6	Features Voting Based on Borda Count	74
3.4.2	Phase 2: Removing Redundant Data Inputs	75
3.4.3	Phase 3: Active-Build Model for Random Forest (ABRF)	78
3.4.3.1	Diagnosis Unit	80
3.4.3.2	Reparation Unit	81
3.4.3.3	Voting Unit	84
3.5	VRF Performance Evaluation Details	84
3.5.1	Internet Traffic Datasets	84
3.5.2	Weka: Data Mining Software in Java	88
3.5.2.1	RapidMiner tool	89

3.5.3	Experiment Stages	89
3.5.4	Experiment Environment	91
3.5.4.1	Hardware Environment	91
3.5.4.2	Software Environment	91
3.5.5	VRF Evaluation Factors	91
3.5.5.1	VRF Evaluation of Phase 1: FSCGRV	92
3.5.5.2	VRF Evaluation Phase 2: Data Reduction	94
3.5.5.3	VRF Evaluation Phase 3: ABRF	95
3.6	Chapter Summary	96
CHAPTER 4 IMPLEMENTATION DETAILS		
4.1	Introduction	97
4.2	Implementation of Proposed FSCGRV Feature Selection	97
4.2.1	Significance Feature Evaluator (SFE)	98
4.2.2	Fast Correlation-Based Features Selection (FCBF)	99
4.2.3	Chi-Squared Feature Selection	100
4.2.4	Gain Ratio (GR)	101
4.2.5	FSCG Based on Ranking and Voting	102
4.3	Removing Redundant Data Inputs	106
4.4	Active Build Model for Random Forest	107
4.4.1	Active and Passive Trees Representation	112
4.5	Chapter Summary	115

CHAPTER 5 RESULTS AND DISCUSSION

5.1	Introduction	116
5.2	Performance Evaluation Metrics for FSCGVR	117
5.2.1	Classification Accuracy Performance for FSCGVR	117
5.2.1.1	Results analysis of classification accuracy based on FSCGRV	125
5.2.2	Processing Time Evaluation for FCSGVR	126
5.2.2.1	Results analysis of processing time based on FSCGRV	137
5.3	Performance Evaluation of Memory Consumption for VRF Based on Data Input Reduction and FSCGRV	138
5.3.1	Results analysis of memory consumption based on DR and FSCGRV	147
5.4	Classification Accuracy and Total Processing Time Performance for VRF Based on ABRF	148
5.4.1	Classification Accuracy Performance for VRF Based on ABRF	149
5.4.2	Total Processing Time for VRF Based on ABRF	154
5.5	The Averaging Performance of Accuracy, Processing Time and Memory Consumption for VRF Over 16 Datasets.	162
5.6	Chapter Summary	166

CHAPTER 6 CONCLUSION AND FUTURE WORK

6.1	Introduction	167
6.2	Conclusion	167
6.3	Suggestions for Future Work	168

REFERENCES

169

LIST OF PUBLICATION

180

©This item is protected by original copyright

LIST OF TABLES

NO.		PAGE
2.1	Port numbers categories.	18
2.2	Signature generated by AutoSig.	22
2.3	Summary of methods used in network traffic classification	58
3.1	Comparison of pre-assessment dataset before and after classification for trees assessment.	81
3.2	Datasets categorization.	86
3.3	Instances involved in source 1 datasets.	87
3.4	Instances involved in datasets source 2 and 3.	88
4.1	Features vector selected by SFE with rank value.	98
4.2	Features vector selected by FCBF with rank value.	99
4.3	Features vector selected by Chi2 with rank value.	100
4.4	Features vector selected by GR with rank value.	102
4.5	The selected features vectors by FSCG with votes and ranking value for Moore datasets.	103
4.6	The votes for each feature over Moore datasets.	103
4.7	The selected features vectors by FSCG with votes and ranking value for source 2 TCP datasets.	104
4.8	The votes for each feature over source2 TCP datasets.	104
4.9	The votes for each feature over data source 3 UDP datasets.	104
4.10	The votes for each feature over source 3 UDP datasets.	105

4.11	The selected features vector in VRF based on FSCG.	105
4.12	Redundant and unique inputs for dataset source1.	106
4.13	Redundant and unique inputs for TCP dataset source2.	107
4.14	Redundant and unique inputs for UDP dataset source3.	107
4.15	Active and passive trees for dataset (Entry 1) with different trees numbers.	108
4.16	Active and passive trees for dataset (Entry 2) with different trees numbers.	108
4.17	Active and passive trees for dataset (Entry 3) with different trees numbers.	108
4.18	Active and passive trees for dataset (Entry 4) with different trees numbers.	109
4.19	Active and passive trees for dataset (Entry 5) with different trees numbers.	109
4.20	Active and passive trees for dataset (Entry 6) with different trees numbers.	109
4.21	Active and passive trees for dataset (Entry 7) with different trees numbers.	109
4.22	Active and passive trees for dataset (Entry 8) with different trees numbers.	110
4.23	Active and passive trees for dataset (Entry 9) with different trees numbers.	110
4.24	Active and passive trees for dataset (Entry 10) with different trees numbers.	110
4.25	Active and passive trees for dataset (TCP G1) with different trees numbers.	111
4.26	Active and passive trees for dataset (TCP G2) with different trees numbers.	111

4.27	Active and passive trees for dataset (TCP G3) with different trees numbers.	111
4.28	Active and passive trees for dataset (UDP G1) with different trees numbers.	111
4.29	Active and passive trees for dataset (UDP G2) with different trees numbers.	112
4.30	Active and passive trees for dataset (UDP G3) with different trees numbers.	112
5.1	Classification accuracy of VRF based on proposed ABRF compared with 5 classifiers for data source 1 (Moore dataset).	151
5.2	Classification accuracy of VRF based on proposed ABRF compared with 5 classifiers for data source 2 (TCP Datasets).	153
5.3	Classification accuracy of VRF based on proposed ABRF compared with 5 classifiers for data source 3 (UDP datasets).	154
5.4	Processing time of VRF based on ABRF compared with 5 classifiers for Data source 2.	158
5.5	Build model time for RF and VRF for data source 1 (Moore datasets).	158
5.6	Processing time of VRF based on ABRF compared with 5 classifiers for Data source 2.	160
5.7	Build model time for RF and VRF for data source 2 (TCP datasets).	160
5.8	Processing time of VRF based on ABRF compared with 5 classifiers for Data source 3.	161
5.9	Build model time for RF and VRF for data source 3 (UDP datasets).	161

LIST OF FIGURES

NO.		PAGE
1.1	Mechanism of supervised and unsupervised machine learning algorithms.	4
1.2	Research methodology.	11
2.1	Structure of chapter two literature review.	14
2.2	Port number of HTTP snapped from Wireshark.	18
2.3	Application detection based on applications signature.	20
2.4	Structure of distributed host based traffic collection platform.	24
2.5	K-means key steps.	29
2.6	One, two and three nearest neighbors.	32
2.7	K-nearest neighbors to the object X.	33
2.8	Probabilistic neural network architecture.	36
2.9	SVM decision boundaries and margins.	37
2.10	Two different decision boundaries in SVM.	37
2.11	Multi class SVM based on one-against-all approach for 4 application.	41
2.12	Multi class SVM based on one-against-one approach for 4 application.	42
2.13	Main steps of expectation maximization.	44
2.14	Hidden Markov Model (HMM) Schematic.	46
2.15	An example of visualization part of C4.5 tree structure for network traffic (from Weka Simulation).	51

2.16	Ensemble classifiers procedure.	54
2.17	Random forests steps and structure.	55
3.1	Vital random forest phases.	60
3.2	Processing sequence of mining data.	61
3.3	FSCGRV Framework.	63
3.4	Pseudo code of fast correlation-based filter.	69
3.5	The voting filter used in FSCGRV.	75
3.6	Labeling steps for redundant inputs in categorization unit.	76
3.7	Flowchart of removing the redundant inputs.	77
3.8	ABRF main units.	78
3.9	The structure of original random forest build model vs. a new ABRF.	79
3.10	Build a new two active trees in ABRF instead the passive trees.	82
3.11	Pseudo code of adding new active tree into forest build model.	83
3.12	Dataset environment.	86
3.13	VRF experiments stages.	90
3.14	VFR framework with regards to evaluation factors.	92

3.15	Classification accuracy evaluation metrics	93
4.1	Passive trees models from different datasets.	113
4.2	Active trees models from different datasets.	114
5.1	The classification accuracy of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the first sixth datasets of Data source 1 (Moore datasets).	119
5.2	The classification accuracy for 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the last fourth datasets of source 1 (Moore datasets).	121
5.3	The classification accuracy for 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the three source 2 datasets (TCP datasets).	122
5.4	The classification accuracy for 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the three source 3 datasets (UDP datasets).	124
5.5	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the first sixth datasets of Data source 1 (Moore datasets).	129
5.6	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the last fourth datasets of Data source 1 (Moore datasets).	131
5.7	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the three Datasets of source2 (TCP datasets).	133
5.8	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the three Datasets of source3 (UDP datasets).	136
5.9	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the first sixth Datasets of source1 (Moore datasets).	140
5.10	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the last fourth Datasets of source1 (Moore datasets).	142

5.11	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the Three Datasets of source2 (TCP datasets).	144
5.12	Processing time of 3 classifiers used 5 features selection techniques compared with proposed FSCGRV for the Three Datasets of source3 (UDP datasets).	146
5.13	The averaging classification accuracy for VRF over 16 datasets compared with 5- classifiers.	163
5.14	The averaging processing time for VRF over 16 datasets compared with 5 classifiers.	164
5.15	The average memory consumption with regard to competing techniques.	165
5.16	The average memory consumption for VRF over 16 datasets compared with 5 classifiers.	165

©This item is protected by original copyright

LIST OF ABBREVIATIONS

ABRF	Active Build-model Random Forest
BFW	Broadband Fixed Wireless
CBMG	Customer Behavior Model Graph
Chi ²	Chi-Squared
CPU	Central Processing Unit
DCCP	Datagram Congestion Control Protocol
DNS	Domain Name System
DoS	Denial of Service
DPI	Deep-Packet Inspection
DR	Data Reduction
DSL	Digital Subscriber Line
EM	Expectation Maximization
FCBF	Fast Correlation-Based Filter
FSCGRV	Fast correlation based filter, Significance feature evaluator, Chi-squared and Gain ratio: and Ranking and Borda Count Voting
FTP	File Transfer Protocol
GR	Gain Ratio
HMM	Hidden Markov Model
HTTP	Hypertext Transfer Protocol
IANA	Internet Assigned Numbers Authority
IBL	Instance-Based Learning
ID3	Iterative Dichotomiser 3
IDS	Intrusion Detection System

IETF	Internet Engineering Task Force
IG	Information Gain
IP	Internet Protocol
IPTs	Inter-Arrival Time Packets
IRC	Internet Relay Chat
ISP	Internet service providers
KNN	K-Nearest Neighbors
ML	Machine Learning
MySQL	My Structured Query Language
NAT	Network Address Translation
NB	Naïve Bayes
NBKE	Naïve Bayes Kernel Estimation
OAA	One-Against-All
OAo	One-Against-One
P2P	Peer-to-Peer
PDF	Probability Density Function
PNN	Probabilistic Neural Network
POP3	Post Office Protocol 3
PS	Packet Size
QoS	Quality of Service
RF	Random Forest
SCTP	Stream Control Transmission Protocol
SFE	Significance Feature Evaluator
SMTP	Simple Mail Transfer Protocol
SSH	Secure Shell

SU	Symmetrical Uncertainty
SVM	Support Vector Machine
TCP	Transmission Control Protocol
UDP	User Datagram Protocol
VoIP	Voice over Internet Protocol
VPN	Virtual Private Network
VRF	Vital Random Forest
WWW	World Wide Web

©This item is protected by original copyright

Pengelasan Trafik Rangkaian Berkesan berdasarkan Vital Random Forest bagi Set Data Dimensi Tinggi

ABSTRAK

Tesis ini mencadangkan serta melaksanakan satu kaedah pengelasan trafik rangkaian yang berkesan berdasarkan *vital random forest* (VRF) baru bagi pemeriksaan data berdimensi tinggi. Kejuruteraan trafik rangkaian merupakan satu daripada teknologi penting yang memaparkan pertumbuhan yang pantas dalam revolusi teknologi seantero dunia. Pengelasan trafik rangkaian memberikan faedah yang boleh dipertimbangkan sebagai suatu wadah kejuruteraan rangkaian yang penting bagi sekuriti rangkaian, reka bentuk rangkaian dan juga bagi pemantauan dan pengurusan rangkaian. Ia memberikan perkhidmatan yang berbeza seperti mengenal pasti aplikasi yang paling banyak menggunakan sumber rangkaian, ia mewakili bahagian teras daripada sistem pengesanan instruksi secara automatik, ia membantu mengesan aplikasi anomali, dan ia juga membantu mengenali aplikasi yang digunakan di seantero dunia bagi tujuan penawaran produk baru. Dalam kata lain, pelbagai cabaran yang dihadapi oleh para jurutera rangkaian dalam usaha mereka mengelaskan trafik. Yang paling biasa adalah penambahan jenis aplikasi dan saiz data trafik yang besar. Oleh itu, berdasarkan kajian lepas, ramai penyelidik berlumba-lumba memperkenalkan kaedah kaedah yang berkesan bagi pengelasan trafik. Keberkesanan pengelasan trafik bergantung pada beberapa faktor penting seperti ketepatan pengelasan, penggunaan memori dan masa pemprosesan. Tesis ini mencadangkan *vital random forest* VRF sebagai pengelasan trafik rangkaian yang berkesan, yang merupakan satu pakej yang memperkenalkan teknik pemilihan penapis baru, pengurangan input data dan model binaan baru bagi kaedah *random forest* baru untuk mengelaskan trafik rangkaian bagi set data yang besar. VRF juga bermatlamat mengurangkan masa pemprosesan, meningkatkan ketepatan pengelasan, serta mengurangkan penggunaan memori. Rangka kerja VRF yang dicadangkan memberikan tiga sumbangan, pertama, reka bentuk *teknik pemilihan penapis baru* berdasarkan penggunaan teknik empat dan dua penapis bagi memilih set penapis yang paling signifikan. Kedua, pengurangan input data yang bertujuan menyingkirkan semua input rekod yang berlebihan dengan pengkatogorian kelas. Ketiga, mereka bentuk model binaan baru bagi *random forest* standard, yang dikenali sebagai ABRF (Active Build model Random Forest). ABRF dibina berdasarkan pokok aktif (pengelas), sebaliknya, pokok pasif didiagnosis, dikeluarkan dan diganti dengan pokok aktif. Keputusan eksperimen adalah berdasarkan beberapa daripada set data tanda aras (data set global). Data ini dikumpul daripada beberapa rangkaian bagi mencapai semua paket yang berkaitan dengan TCP, UDP dan IP dalam kedua-dua arah. Keputusan menunjukkan penambahbaikan yang ketara dari segi faktor yang signifikan, iaitu ketepatan pengelasan, masa pemprosesan dan penggunaan memori. Secara purata, keputusan VRF berhubung dengan faktor ini dijalankan berdasarkan 16 set data tanda aras, ketepatan keseluruhan VRF meningkat sebanyak 6% berbanding dengan hutan rawak asal untuk mencapai 99.58%, masa pemprosesan telah menurun dengan perbezaan 62% manakala VRF menggunakan hanya 38% daripada jumlah purata masa dan purata penggunaan memori dikurangkan sebanyak 40%. Selanjutnya, VRF telah disahkan melalui perbandingan keputusan daripada faktor yang dinyatakan di atas dengan penyelidikan terdahulu.

An Efficient Network Traffic Classification Based on Vital Random Forest for High Dimensional Dataset

ABSTRACT

This thesis proposes and implements an efficient network traffic classification method based on a new vital random forest for high dimensional data. Network traffic engineering is one of the most important technologies that have witnessed a rapid growth in the revolution of worldwide technologies. Network traffic classification has added considerable interest as an important network engineering tool for network security, network design, as well as network monitoring and management. It can introduce different services such as identifying the applications which are most consuming for network resources, it represents the core part of automated intrusion detection systems, it helps to detect anomaly applications and it helps to know the widely-used applications for the intention of offering new products. On the other hand, several challenges faced by network engineers on their course to classify traffic. The most common of which are increasing application types and the huge size of data traffics. Therefore, many researchers have been competing in literature to introduce an efficient method for traffic classification. The efficiency is dependent on important factors such as classification accuracy, memory consumption and processing time. This thesis presents a *Vital Random Forest (VRF)* as efficient network traffic classification which is a one package that introduces a new features-selection technique, data inputs reduction and a new build model for original random forest method to classify network traffic for huge datasets. VRF aims to reduce processing time, increase classification accuracy and decrease memory consumption. The proposed framework of VRF consists of three contributions; first one is design of a *new features-selection technique* based on adopting four techniques and two filters for selecting most significant features set. Second is a data input reduction aiming to remove all redundant record inputs with class categorization. Third is to design a new build model for standard random forest called *Active Build model Random Forest (ABRF)*. ABRF is built based on active trees (classifiers) while the passive trees are diagnosed, omitted and replaced with active trees. The results from the experiments conducted are based on several benchmark datasets (global dataset). These data were collected from the edge of a network to access all packets associated with a TCP, UDP and IP connections in both directions. The results exhibit noticeable improvement in terms of three significant factors, namely; classification accuracy, processing time and memory consumption. The averaging results of VRF with regard to these factors were conducted based on 16 benchmark datasets where the classification accuracy of VRF is increased by 6% compared with original random forest to reach 99.58%, the processing time was decreased down with difference 62% while VRF consumed only 38% from the total average time and the averaging memory consumption is reduced by 40%. Furthermore, VRF has been validated by comparing the results for the above-mentioned factors with previous works.

CHAPTER 1

INTRODUCTION

1.1 Introduction

Network traffic engineering is one of the most important technologies that have witnessed a rapid growth in the revolution of global technologies. Network traffic identification and classification have recently gained considerable interest as an important network engineering tool for network security, network design, as well as network monitoring and management (Dainotti, Pescapé, & Claffy, 2012).

Network traffic classification is a process that categorizes network traffic according to various parameters (e.g. port number, arrival time, type of protocol, packet length, etc) into a number of application classes such as (P2P, WWW, Mail, etc.). This method introduces multi-beneficial solutions in different avenues, such as network security, network management, network measurement and quality-of-service (QoS) (Callado et al., 2009). Many of network engineers have had started to inspect and analyze network traffic but they faced several novel challenges the most significant of which are the huge amount of classified data traffic and the variety of applications. As a result, numerous studies have been proposed to face the above mentioned challenges and they started competing in terms of different factors such as accuracy of classification, memory consumption, CPU consumption and time consumption of processing (Dainotti, Pescapé, & Sansone, 2011; Nguyen & Armitage, 2008). The studies introduced various solutions that involved two major types: First, representing traditional solutions that include well-known port number, deep packet inspection and behavior-based approach. Second is Machine Learning (ML): there was a considerable interest in several disciplines such as

philosophical, logical and conceptual issues, and then research interest shifted to computational and algorithmic aspects of ML that is driven mainly by practical application. That's why many studies were rerouted towards ML. This thesis focuses on supervised machine learning method and its contribution in the area of network traffic classification.

1.2 Background

This section provides a general background to the work presented in this thesis. It briefly introduces the principal technologies referenced throughout this thesis as network traffic engineering, Machine Learning and Random Forest.

1.2.1 Network Traffic Engineering

First and foremost, to understand the network traffic engineering we need to define the network engineering, it's a method of manipulating your network to suit your traffic. Network traffic engineering is defined as that aspect of network engineering of dealing with the issue of performance evaluation and performance optimization of operational IP networks (Callado et al., 2009). Traffic Engineering encompasses the application of technology and scientific principles to the measurement, characterization, modeling, and control of network traffic. The enhancement of the network is achieved at both levels: traffic and resources, with regard to different performance requirements and utilizing network resources economically and reliably. The performance is measured in terms of delay, jitter, packet loss and throughput. The main purpose of network traffic engineering is to facilitate reliable network operations. This is accomplished using policies to keep network survivability and minimizing the vulnerability of the network to service outages