

COMPARISON OF SPAM FILTERING MECHANISM

LESTA ANAK ENKAMAT



UNIVERSITI TEKNIKAL MALAYSIA MELAKA

BORANG PENGESAHAN STATUS TESIS *

JUDUL: COMPARISON OF SPAM FILTERING MECHANISM

SESI PENGAJIAN: 2010/2011

Saya LESTA ANAK ENKAMAT

(HURUF BESAR)

mengaku membenarkan tesis (PSM/ Sarjana/ Doktor Falsafah) ini disimpan di Perpustakaan Teknologi Maklumat dan Komunikasi dengan syarat-syarat kegunaan seperti berikut:

1. Tesis dan projek adalah hakmilik Universiti Teknikal Malaysia Melaka.
2. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan untuk tujuan pengajian sahaja.
3. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan tesis ini sebagai bahan pertukaran antara institusi pengajian tinggi.
4. ** Sila tandakan (/)

 SULIT

(Mengandungi maklumat yang berdarjah keselamatan atau kepentingan Malaysia seperti yang termaktub di dalam AKTA RAHSIA RASMI 1972)

 TERHAD

(Mengandungi maklumat TERHAD yang telah ditentukan oleh organisasi/ badan di mana penyelidikan dijalankan)

 / TIDAK TERHAD



(TANDATANGAN PENULIS)

Alamat Tetap: Rumah Usek,
Sungai Kriok, 98200 Niah,
Sarawak.

Tarikh: 12 Julai 2011



(TANDATANGAN PENYELIA)

En. Mohammad Radzi bin Motsidi

Tarikh: 12 Julai 2011

CATATAN: *Tesis dimaksudkan sebagai Laporan Akhir Projek Sarjana Muda (PSM)

**Jika tesis ini SULIT atau TERHAD, sila lampirkan surat daripada pihak berkuasa.

COMPARISON OF SPAM FILTERING MECHANISM

LESTA ANAK ENKAMAT

**This report is submitted in partial fulfillment of the requirements for the
Bachelor of Science Computer (Computer Networking)**

**FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY
UNIVERSITI TEKNIKAL MALAYSIA MELAKA**

2011

DECLARATION

I hereby declare that this project report entitled
COMPARISON OF SPAM FILTERING MECHANISM

is written by me and is my own effort and that no part has been plagiarized
without citations.

STUDENT :  Date: 12 JULAI 2011
(LESTA ANAK ENKAMAT)

SUPERVISOR :  Date: 12 JULAI 2011
(EN. MOHAMMAD RADZI BIN MOTSIDI)

DEDICATION

To my beloved parents, your care, your love and your support give me strength.

To my friends, it is for your support and motivation.

To my lecturer, for the guide until the completion of this project, it gives me the challenge to be a better person.

ACKNOWLEDGEMENTS

During the completion of this project, I gained a lot of new knowledge especially in email system. I also learned a new things in life, be strong whenever face a problem and do not give up because of failed many times before. I would like to thank all people that involved with this project.

To my supervisor Mr. Mohammad Radzi bin Motsidi, I would like to thank him for assisting me in completion of this project. He helped me a lot by giving me ideas and suggestions on how to run this whole project. I really appreciate it.

Apart from that, I also want to express my sincere thanks to my all my friends who gave me the support and motivation during the completion of this project.

Lastly, I would like to express my deepest thanks to my family especially my mum who has give me support in completion of this project.

ABSTRACT

Email spam or unsolicited bulk email (UBE) that usually sent to many people either through email. Spam usually wasting our time by deleting it one by one in our mail box, and also eats a lot of network bandwidth. Nowadays, there are lots of organizations and individuals take a step to fight spam with variety of techniques. This research paper is focus on two different types of spam filtering software which are SpamAssassin and ThunderBayes. A testing done to make comparison on both of this spam filtering software based on how this program filter and classified email as spam. The results of test show which program is effective to be used as spam filter software on mail server.

ABSTRAK

Email spam ataupun *unsolicited bulk email (UBE)* biasanya dihantar ke ramai orang sekaligus melalui emel. Spam selalunya membazir masa kita kerana kita perlu membuang emel tersebut satu per satu di dalam Inbox emel kita. Selain itu, spam juga memakan banyak *bandwidth* rangkaian. Pada masa kini, kebanyakan organisasi dan juga individu mengambil langkah untuk melawan spam dengan pelbagai teknik. Kajian ini tertumpu pada dua jenis program yang digunakan untuk menapis emel spam, iaitu *SpamAssassin* dan *ThunderBayes*. Ujian dilakukan untuk membandingkan bagaimana kedua-dua program ini menapis dan mengklasifikasikan emel tersebut adalah spam. Dari keputusan ujian tersebut, dapatlah diketahui mana satu program yang lebih efektif untuk menapis emel spam.

TABLE OF CONTENTS

CHAPTER	SUBJECT	PAGE
	DECLARATION	ii
	DEDICATION	iii
	ACKNOWLEDGEMENTS	iv
	ABSTRACT	v
	ABSTRAK	vi
	TABLE OF CONTENTS	vii
	LIST OF TABLES	xii
	LIST OF FIGURES	xiii
	LIST OF ABBREVIATIONS	xvi
	LIST OF APPENDICES	xvii
CHAPTER I	INTRODUCTION	
	1.1 Project Background	1
	1.2 Problem Statements	2
	1.3 Objective	2
	1.4 Scope	3
	1.5 Project Significance	3
	1.6 Expected Output	4
	1.7 Conclusion	4
CHAPTER II	LITERATURE REVIEW AND PROJECT METHODOLOGY	

2.1	Introduction	5
2.2	Literature Review	5
2.2.1	Domain	6
2.2.2	Keyword	6
2.2.2.1	Spam	6
2.2.2.2	Client-server	7
2.2.2.3	Email server	7
2.2.2.4	Email client	7
2.2.2.5	Email filtering	8
2.2.2.6	Spam filter	8
2.2.3	Previous Research	8
2.2.3.1	Related Research	8
2.2.3.1.1	Revisit Bayesian	8
	Approaches for Spam Detection	
2.2.3.1.2	A Case-Based Approach to Spam Filtering That Can Track Concept Drift	9
2.2.3.1.3	A Classification Method for Spam E-mail by Self-Organizing Map and Automatically Defined Groups	9
2.2.3.2	Spam Filtering Techniques	10
2.2.3.2.1	Rule-based scoring techniques	10
2.2.3.2.2	Collaborative spam filtering techniques	11
2.2.3.2.3	Bayesian filtering techniques	11
2.2.3.3	Comparison of spam filtering tool	12
2.2.3.4	Email System – How it works?	13
2.2.3.5	Basic parts of Email system	13
2.2.3.6	Protocols for Email	15
2.2.3.7	How SpamAssassin Works?	16
2.2.3.8	How ThunderBayes Works?	17
2.3	Proposed Solution	17

2.3.1	Project Methodology	17
2.4	Project Schedule and Milestones	19
2.4.1	Project Schedule	19
2.4.2	Gantt Chart	22
2.5	Conclusion	22

CHAPTER III ANALYSIS

3.1	Introduction	23
3.2	Analysis of Current Email System	23
3.3	Requirement Analysis of New Email System with Spam Filtering Method	24
3.3.1	Network Architecture	25
3.3.2	Logical and Physical Design	26
3.3.3	Quality of Data	27
3.4	Conclusion	28

CHAPTER IV DESIGN

4.1	Introduction	29
4.2	Possible Scenarios	29
4.2.1	Scenario One	30
4.2.1.1	Email system without spam filtering installed inside	30
4.2.1.2	Software Requirement	31
4.2.1.3	Flowchart (Without spam filter)	32
4.2.2	Scenarion Two	34
4.2.2.1	Email system with SpamAssassin installed	34
4.2.2.2	Software Requirement	35
4.2.2.3	Flowchart (with SpamAssassin)	36
4.2.3	Scenario Three	38
4.2.3.1	Email system with ThunderBayes	38

	installed	
	4.2.3.2 Software Requirement	39
	4.2.3.2 Flowchart (with ThunderBayes)	40
4.3	Conclusion	41

CHAPTER V IMPLEMENTATION

5.1	Introduction	42
5.2	Network Configuration Management	43
	5.2.1 Configuration Environment Setup	43
	5.2.1.1 BIND DNS Server	43
	5.2.1.2 Postfix Mail Server	47
	5.2.1.3 Dovecot IMAP/POP3 Server	49
	5.2.1.4 Procmail Mail Filter (SpamAssassin)	50
	5.2.1.5 SpamAssassin	50
	5.2.1.6 ThunderBayes	51
5.3	Hardware Configuration Management	53
	5.3.1 Hardware Setup	53
5.4	Development Status	53
5.5	Conclusion	54

CHAPTER VI TESTING

6.1	Introduction	55
6.2	Test Plan	55
	6.2.1 Test Organization	56
	6.2.2 Test Environment	56
	6.2.3 Test Schedule	56
6.3	Test Strategy	57
	6.3.1 Classes of tests	58
6.4	Test Design	58
	6.4.1 Test Description	59
	6.4.2 Test Data	59

6.5	Test Result and Analysis	69
6.6	Conclusion	74

CHAPTER VII PROJECT CONCLUSION

7.1	Observation on Weakness and Strength	75
7.1.1	Strengths	75
7.1.2	Weaknesses	76
7.2	Propositions for Improvement	76
7.3	Contribution	77
7.4	Conclusion	78
	REFERENCES	79
	BIBLIOGRAPHY	80
	APPENDIX A	81
	APPENDIX B	82
	APPENDIX C	88

LIST OF TABLES

TABLE	TITLE	PAGE
2.1	Comparison of Spam Filtering Software	12
2.2	Project Schedules and Milestone	20
4.1	Table of Software Requirement without Spam Filter	31
4.2	Table of Software Requirement with SpamAssassin.	35
4.3	Table of Software Requirement with ThunderBayes	39
5.1	Hardware configuration inside virtual machine	53
6.1	Table of Configuration of Virtual Machine inside VMware	56
6.2	Table of Test Schedule for Each Test	57
6.3	Table of Classes of Tests	58
6.4	Table of Test Description	59
6.5	Table of DNS server detail for each scenario	62
6.6	Summary of Test Result for SpamAssassin and ThunderBayes	73

LIST OF FIGURES

FIGURE	TITLE	PAGE
2.1	Schematic representation of SpamAssassin	16
3.1	Email system architecture without spam filtering method	23
3.2	Email system architecture with spam filtering method	24
3.3	Client-Server architecture	25
3.4	Logical Network Diagram for email server with SpamAssassin	26
3.5	Logical Network Diagram for email server with ThunderBayes	26
3.6	Physical Network Diagram	27
4.1	Email server without spam filter	30
4.2	Flow Chart for Server Side (no filter)	32
4.3	Flow Chart for Client Side (no filter)	33
4.4	Email System with SpamAssassin	34
4.5	Flow Chart for Server Side (SpamAssassin)	36
4.6	Flow Chart for Client Side (SpamAssassin)	37
4.7	Email system with ThunderBayes	38
4.8	Flow Chart for Server Side (ThunderBayes)	40
4.9	Flow Chart for Client Side (ThunderBayes)	41
5.1	Main configuration setting for DNS server (SpamAssassin)	44
5.2	Forward Lookup Zone setting for DNS server (SpamAssassin)	44

5.3	Reverse Lookup Zone setting for DNS server (SpamAssassin)	45
5.4	Main configuration setting for DNS server (ThunderBayes)	45
5.5	Forward Lookup Zone configuration for DNS server (ThunderBayes)	46
5.6	Reverse Lookup Zone for DNS server (ThunderBayes)	46
5.7	Main.cf configurations for Postfix mail server (SpamAssassin)	47
5.8	Master.cf configuration for Postfix mail server (SpamAssassin)	48
5.9	Main.cf configurations for Postfix mail server (ThunderBayes)	48
5.10	Master.cf configurations for Postfix mail server (ThunderBayes)	49
5.11	Dovecot configuration for IMAP/POP3 server	49
5.12	Procmail configuration for SpamAssassin	50
5.13	Local.cf configuration for SpamAssassin	51
5.14	ThunderBayes Options for Mozilla Thunderbird.	52
5.15	Training page for ThunderBayes	52
6.1	Ping from Server to Client (SpamAssassin)	60
6.2	Ping from Client to Server (SpamAssassin)	60
6.3	Ping from Server to Client (ThunderBayes)	61
6.4	Ping from Client to Server (ThunderBayes)	61
6.5	DNS testing with 'dig' command for mail.finalproject.com (SpamAssassin)	62
6.6	DNS testing with 'dig' command for mail.bayesproject.com (Thunderbayes)	63
6.7	Telnet to mail.finalproject.com (SpamAssassin)	64

6.8	Telnet to mail.bayesproject.com (ThunderBayes)	64
6.9	A ham email sample sent from one client to another client	65
6.10	Client received the ham email sample in Inbox folder	65
6.11	A spam email sample sent from one client to another client	66
6.12	Client received the spam email sample in Junk folder	66
6.13	A ham email sample sent from one client to another client	67
6.14	Client received the ham email sample in Inbox folder	67
6.15	A spam email sample sent from one client to another client	68
6.16	Client received the spam email sample in Junk folder	68
6.17	Fifty samples of ham email in Inbox folder	69
6.18	Forty-seven samples spam email in Junk folder	70
6.19	Three samples spam email in Inbox folder	70
6.20	Forty-seven samples of ham email in Inbox folder	71
6.21	Three samples ham email in Junk folder	72
6.22	Forty-nine samples of spam email in Junk folder	72
6.23	Spam email that classify as ham in Inbox.	73

LIST OF ABBREVIATIONS

SMTP	Simple Mail Transfer Protocol
MTA	Mail Transfer Agent
MUA	Mail User Agent
POP3	Post Office Protocol 3
IMAP	Internet Message Access Protocol
IP	Internet Protocol
DNS	Domain Name System
FQDN	Fully Qualified Domain Name
UBE	Unsolicited Bulk Email
UCE	Unsolicited Commercial Email
TCP	Transmission Control Protocol
LAN	Local Area Network
TCP/IP	Transmission Control Protocol/Internet Protocol
SOM	Self-Organizing Map
ADG	Automatically Defined Groups
ASSP	Anti-Spam SMTP Proxy
PSFL	Python Software Foundation License
GNU	General Public License
MIME	Multipurpose Internet Mail Extensions
HTML	Hyper Text Markup Language
OS	Operating System

LIST OF APPENDICES

APPENDIX	TITLE
Appendix A	Gantt Chart
Appendix B	Project Proposal
Appendix C	Log Book

CHAPTER I

INTRODUCTION

This chapter is the early reviews for this project, which will explain about the project that will be developed, that is spam filtering mechanism in email system. Sub-chapters that will be discussed in this chapter including project background, problem that related to this project, objectives of this project, scopes of the project, project significance, expected output of this project and the conclusion of this chapter.

1.1 Project Background

Most spam emails on the Internet today are advertisements from individuals and the occasional small business that looking for a way to make a fast money. E-mail spam, which also known as junk e-mail or unsolicited bulk e-mail (UBE) is a subset of spam that involves nearly identical messages sent to numerous recipients by e-mail. Definitions of spam usually include the aspects that e-mail is unsolicited and sent in bulk. Usually when user log in to their email, they can see few spam e-mail go through to their inbox. This spam sometime annoying because when it comes with large amount of spam email, it will take time to remove all the junk mail and also it use a lot of storage space on the email. Some spam emails also can contain malware. For regular users which don't have knowledge about malware, they just simply open the spam email that may contain malware. The malware can affect the user's computer and worst case which can happen is the malware will spread itself to

the entire network of the company or organization if the computer is connected to a company networks.

To overcome this problem, the spam filtering mechanism is important because it can recognize the spam and block the spam e-mail or move the spam e-mail to the Junk folder in client email.

For this project, spam filtering implementation will use two types of spam filtering tools, which are SpamAssassin and ThunderBayes. These two applications will be installed and configured alongside with email system so that it can used to filter spam emails. Testing will be done to compare the usage of the software and how the software work to identify email that classified as spam or ham.

1.2 Problem Statement

Spam emails wasting the Internet's two most precious resources, which are the bandwidth of long-distance communications links and the time of network administrator who keep the Internet working from day to day. Spam emails also wasting the time of computer users around the world, which they need to delete those spam emails. Furthermore, in order to deliver the emails, persons who send spam emails are increase which can lead to fraud and computer abuse.

1.3 Objective

There are few main objectives to be established during completion of this project. The objectives have been summarized as below:

- i. To setup and configure mail server.
- ii. To implement the spam filtering tool into mail server that will detect and filter spam i

- iii. To compare the implemented spam filtering mechanism software. The spam filtering tools that will be used is SpamAssassin and ThunderBayes.

1.4 Scope

The scope of this project is to implement and compare spam filtering tool which can be used to detect spam in email. The software that will be used to filter spam emails are SpamAssassin and ThunderBayes. The spam filtering tool will be implemented in Linux operating system environment. This spam filtering mechanism is recommended for any organization or small companies that provide email service to their users or workers.

1.5 Project Significance

This project will implement email system with spam filtering tool. This spam filtering tool will detect spam emails that sent to user's email box. For a large organization or perhaps a small company, an email system is an essential way to deliver latest information to employees. When there's update or new task that need to be done in the organization, administrator or whoever in charged with it can inform the workers inside the organization using emails.

For a large company which have a lot of branch, Internet connection needed in order to connect the email system. When it connected to Internet, the risk of being attack with spam emails is high. So, spam email mechanism needed to handle those spam emails. Because some spam email can contain malware which can lead to fraud, it can affect the security of the large organization. Malware such as spyware can stole information which stored inside the organization's server.

When the spam filtering tool implemented inside an email server, it can help large organization company filters spam emails.

1.6 Expected Output

This project will implement email system with spam filtering tool. This spam filtering tool will able to detect and filter spam or ham emails that incoming to user's mail box. Two types of spam filtering tool will be used and comparison should be made based of how the software filters the incoming emails. For SpamAssassin, it will implement on server side while ThunderBayes will implement in client side. Testing will be done in order to compare these two types of spam filtering tool.

1.7 Conclusion

In conclusion, this chapter describes the background of this project that will be developed. This project will used two types of open source spam filtering tool to detect and filter incoming spam or ham. The spam filtering tool will be implemented alongside with email server inside a Linux based operating system.

Next chapter will be discussing about the literature review and project methodology that will be used to implement this project.

CHAPTER II

LITERATURE REVIEW AND PROJECT METHODOLOGY

2.1 Introduction

The aim of this project is to review previous research that related to this project. The most relevant methodology that will be used to develop this project also will be identified in this chapter. Several researches are being conducted to gain more understanding and ideas in order to complete this project.

The methodology that will be used through completion of this project is Top-Down approach.

2.2 Literature Review

According to the objectives and scopes of this project, this literature review will give a summary about previous research about problem with spam emails and how to handle the problem. This will give reader what knowledge and ideas have been done that related to this project and also state the strengths and weaknesses of the project.