

**LAYERED BOTNET DETECTION FRAMEWORK BASED ON SIGNAL
PROCESSING AND DISCRETE TIME ANALYSIS**

By

LOAI KAYED HASSAN BANI MELHIM

UNIVERSITI SINSE MALAYSIA

September 2012

**Layered Botnet Detection Framework Based on Signal Processing
and Discrete Time Analysis**

By

LOAI KAYED HASSAN BANI MELHIM

**Thesis submitted in fulfillment of the
Requirements for the degree of
Doctor of philosophy**

September 2012

ACKNOWLEDGMENTS

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

وَعَدَ اللَّهُ الَّذِينَ آمَنُوا مِنْكُمْ وَعَمِلُوا الصَّالِحَاتِ لَيَسْتَخْلِفَنَّهُمْ فِي الْأَرْضِ كَمَا اسْتَخْلَفَ الَّذِينَ مِنْ قَبْلِهِمْ وَلَيُمَكِّنَنَّ لَهُمْ دِينَهُمُ الَّذِي ارْتَضَى لَهُمْ وَلَيُبَدِّلَنَّهُمْ مِنْ بَعْدِ خَوْفِهِمْ أَمْنًا يَعْبُدُونَنِي لَا يُشْرِكُونَ بِي شَيْئًا وَمَنْ كَفَرَ بَعْدَ ذَلِكَ فَأُولَئِكَ هُمُ الْفَاسِقُونَ (55) النور

All Praise go to almighty Allah, we thank him, seek his guidance and forgiveness for giving me the patience and health to finish this research.

There have been many hands helping me in my work, and some people stand out in my mind. The first idea related to this thesis came from my friend and my previous supervisor, Dr Ahmad Manasreh. Who helped, encouraged and support me. If the meaning of pursuing a Ph.D. is more than a degree, all of the lessons I learned from this journey has permeated my whole life. I feel profoundly enriched and fortunate to have met all of those whom I interacted with.

I would like to express my deepest gratitude to my main supervisor, Assoc. Professor Wan TatChee for his valuable insights, comments, patient guidance and continual encouragement and to my co-supervisor Professor Sureswaran Ramadass Director of NAV6 Center of Excellence for his invaluable guidance through the course of my research. Without their knowledge and assistance, this thesis would not have been successful.

Special thanks to National Advanced IPv6 Center of Excellence (NAv6), and to all NAv6 staff and to the whole NSST team especially to Hairul Nizar, Othman, Kunalan dava Rajoo. Without your time, patience and resources I would never have been able to accomplish this research.

Thanks to my friends Ammar al-momani, Khaleel Shaqfa, Eisa, Sami, mohammad, and all the staff in Al-Zelfi colledge and to all my colleagues for making my research days happy and memorable. I also would like to express my gratitude to Eng.Muyad bani melhim for his constant support in Matlab programming as I have pursued my research.

I would like to express my deep thanks to those who are close to my heart; my brothers, my sisters, and all of my family members, thank you so much for your love and support.

Above all, I would like to express my love and gratitude to my wife, Rasha, and my dearest son Ahmad, for their understanding and endless love through the duration of my studies.

I dedicate this thesis sincerely to my father, and my mother who guide me and helped me to be what I am. Without them, I would not have reached what I have reached today. For their continuous support and guidance I owe much to their patience, and their honest advices. It was their many encouragements that kept me going when I had doubts about myself.

تمت بحمد الله

Loai Kayed Bani Melhim

Penang, Malaysia. March 2013.

TABLE OF CONTENTS

Acknowledgements	ii
Table of Contents	iv
List of Tables	x
List of Figures	xiv
List of Abbreviations	xxiii
List of Appendices	xxv
Abstrak	xxvi
Abstract	xxviii

CHAPTER 1: INTRODUCTION

1.1	The Scale of Botnet Problem	2
1.2	Research Motivation	3
1.3	Problem Statement	4
1.4	Research Objectives	5
1.5	Thesis Contribution	6
1.6	Thesis Scope	7
1.7	Research Framework	8
1.8	Thesis Outline	9

CHAPTER 2: LITERATURE REVIEW

2.1	Introduction	10
2.2	Background of Botnet	11
2.2.1	Definitions Related to Botnet	11
2.2.2	Current and Expected Future Structure of Botnet	12
2.2.3	How Bots Work	17

2.2.4	Botnet evolution	18
2.2.5	Botnet Potential Benefits	21
2.2.6	Botnet Infection Mechanisms	21
2.3	Botnet Life Cycle & Detection Systems	21
2.4	Botnet Detection Methods	24
2.4.1	Signature-based Detection	24
2.4.2	Anomaly-based Detection	26
2.4.3	DNS-based Detection	29
2.4.4	Mining-based Detection	30
2.5	Periodic and Non Periodic signals	35
2.5.1	Periodic Signals	35
2.5.2	Aperiodic (Non Periodic) Signals	36
2.6	Periodic Behavior Detection	36
2.6.1	Power Spectral Density (PSD)	37
2.6.2	Periodograms	37
2.6.3	The Circular Autocorrelation Function	40
2.6.4	Periodogram Vs Autocorrelation Function	42
2.7	Network Traffic Capturing Tools	44
2.8	Traffic Identifier	45
2.9	Summary	49

CHAPTER 3: DESIGN AND METHODOLOGY

3.1	Introduction	52
3.2	Layered Botnet Detection Framework Architecture	52
3.2.1	Botnet Detection Framework Overview	54
3.2.2	Layered Botnet Detection Framework (LBDF)	55
3.3	Layer One: Network Traffic Capturing	56

3.4	Layer Two: Traffic Identifier	57
	3.4.1 Network Traffic Decomposition into SNAK and Data Planes	58
	3.4.2 Generating Discrete Time Sequences $x[n]$	60
	3.4.3 SNAK and data Similar Behavior	61
3.5	Layer Three: Periodic Behavior Detection	65
	3.5.1 Time Domain to Frequency Domain	66
	3.5.2 Periodic Behavior Detection	66
	3.5.2.1 Periodicity Detection	67
	3.5.2.2 Periodicity Validation	68
3.6	Layer Four: Malicious Activity Detector	71
	3.6.1 Inbound Scan Detection (ISD)	72
	3.6.1.1 Evaluation of Scan Attempt Weights (Whs & Wls)	73
	3.6.1.2 Evaluation of LBDF Scan-warning threshold (Sth)	73
	3.6.2 Outbound Scan Detection (OSD)	76
3.7	Layer Five: Analyzer	79
3.8	Chapter Summary	80
 CHAPTER 4: IMPLEMENTATION		
4.1	Introduction	81
4.2	Layered Botnet Detection Framework Architecture	81
4.3	Layer One: Network Traffic Capturing	84
4.4	Layer Two: Traffic Identifier	84
	4.4.1 Network Traffic Decomposition into SNAK and Data Planes	85
	4.4.2 Generating Discrete Time Sequences $x[n]$	86
	4.4.3 SNAK and data Similar Behavior	88

4.5	Layer Three: Periodic Behavior Detection	90
4.5.1	Time Domain to Frequency Domain Conversion	90
4.5.2	Periodic Behavior Detection	91
4.5.2.1	Periodogram of X[k]	91
4.5.2.2	Candidate Peaks Discovery	92
4.5.2.3	Specifying candidate peaks locations	93
4.5.2.4	Circular Autocorrelation	94
4.5.2.5	Candidate Peaks Verification	95
4.5.2.6	Identification of the Fundamental Frequency	100
4.5.2.7	Periodic Behavior verification	101
4.6	Layer Four: Malicious Activity Detector	103
4.6.1	Inbound Scan Detection (ISD)	103
4.6.1.1	Evaluation of Scan Attempt Weights (Whs & Wls)	103
4.6.1.2	Evaluation of LBDF Scan-warning threshold (Sth)	103
4.6.2	Outbound Scan Detection (OSD)	107
4.7	Layer Five: Analyzer	110
4.8	Chapter Summary	112

CHAPTER 5: EXPERIMENTAL ENVIROMENT

5.1	Introduction	113
5.2	Experimental Environmental for LBDF	112
5.3	Verification of Layer One: Network Traffic Capturing	114
5.4	Verification of Layer Two: Traffic Identifier	114
5.4.1	Network Traffic Decomposition into SNAK and Data Planes	114
5.4.2	SNAK and data Similar Behavior	116
5.5	Verification of Layer Three: Periodic Behavior Detection	117

5.6	Verification Layer Four: Malicious Activity Detector	118
5.6.1	Inbound Scan Detection (ISD)	118
	5.6.1.1 Experiments to Evaluate LBDF Scan-warning threshold (Sth)	118
	5.6.1.2 Experiments to Verify Inbound Scan Detection (ISD)	120
5.6.2	Outbound Scan Detection (OSD)	120
5.7	Experiments to Verify LBDF	123
5.7.1	Experiments to Distinguish between Normal and Malicious Files	123
5.7.2	Experiments to Validate LBDF Performance	125
5.8	LBDF Detection Evaluation Methods	127
5.9	Chapter Summary	128

CHAPTER 6: RESULTS AND DISCUSSION

6.1	Introduction	130
6.2	Layered Botnet Detection Framework	130
6.3	Layer One: Network Traffic Capturing	131
6.4	Layer Two: Traffic Identifier	131
	6.4.1 Network Traffic Decomposition into SNAK and Data Planes	132
	6.4.2 Time Variation of Cross-Correlation	133
6.5	Layer Three: Periodic Behavior Detection	135
	6.5.1 Periodic Signals with Single Frequency	135
	6.5.2 Solar Radiation Dataset	141
6.6	Layer Four: Malicious Activity Detector	144
	6.6.1 Inbound Scan Detection (ISD)	144

6.6.1.1	Evaluation of LBDF Threshold Scan-warning threshold (Sth)	144
6.6.1.2	Inbound Scan Detection Results	147
6.6.2	Outbound Scan Detection Results	148
6.7	LBDF Results	149
6.7.1	LBDF Bot Detection Results	150
6.7.2	6.7.2Evaluation of LBDF vs other Detection Methods	167
6.7	Chapter Summary	169
CHAPTER 7: CONCLUSION AND FUTURE WORK		
7.1	Conclusion	170
7.2	Contributions	170
7.2	Future Work	172
REFERENCES		173
APPENDIX A: Layer Two: Traffic Identifier Results		180
APPENDIX B: LBDF Evaluation Results		193
LIST OF PUBLICATIONS		210

LIST OF TABLES

		Page
Table 1.1	State of Botnet at the end Of 2010	3
Table 2.1	C&C topologies and basic properties	15
Table 2.2	Comparison of Periodogram and ACF for periodicity Detection	43
Table 2.3	Botnet Detection Methods Summarization	50
Table 3.1	High severity hs ports	77
Table 4.1	Expected slope values of each ACF segments.	97
Table 5.1	Description of NAv6 dataset Files	115
Table 5.2	Description of DARPA 1999 dataset.	116
Table 5.3	Normal traffic resulted by Set Experiment I	119
Table 5.4	Inbound Scan traffic resulted by Set Experiment II Part1	120
Table 5.5	Outbound Scan traffic resulted by Set Experiment II Part2	122
Table 5.6	Files that used to validate LBDF and their sizes and packet counts.	125
Table 5.7	The used Dataset15 , consist of 15 instances of normal and different malicious traffic.	126
Table 5.8	Dataset15 description.	126
Table 5.9	Dataset15 description.	127
Table 6.1	the k values of first and second segment S1 and S2 and there corresponding values of ACF	139
Table 6.2	The list of candidate peaks for the solar radiation dataset at $\alpha = 0.01$	142
Table 6.3	The k values of P2-segments, S1 and S2 and there corresponding values of ACF.	143
Table 6.4	LBDF layer three testing results by applying solar dataset and sinusoidal signal.	144
Table 6.5	Example of daily calculated Sth.	145

Table 6.6	the daily average Sth for each experiment period.	145
Table 6.7	The final calculated average Sth for the whole period.	145
Table 6.8	The new labeling scheme based on S value, the range of S is selected based on the observed S values in our experiments.	146
Table 6.9	Set Experiment II Part1 files, traffic reduction resulted by Layer 2 algorithms	147
Table 6.10	The result of inbound scan detection on the RegularT1 File.	148
Table 6.11	The Outbound detection results of file OnehostT1	149
Table 6.12	The output of the malicious activity detector	149
Table 6.13	<p> SNAK is used to reduce the amount of traffic to be analyzed. For Zeus file the reduction was with a percentage of 24%. </p>	150
Table 6.14	The list of candidate peaks resulted by the periodogram.	153
Table 6.15	The verified peaks by ACF, only one peak is considered and the rest are discarded.	154
Table 6.16	The Outbound scanning result of the monitored host.	154
Table 6.17	The verified peaks by ACF, only one peak is considered, and the rest are discarded.	156
Table 6.18	List of candidate peaks after been verified by ACF.	156
Table 6.19	The list of candidate peaks resulted by the periodogram.	158
Table 6.20	The Outbound scanning result of the monitored host.	158
Table 6.21	<p> SNAK is used to reduce the amount of traffic to be analyzed. For Scanner host file the reduction was with a percentage of 50%. </p>	159
Table 6.22	The Outbound scanning result of the monitored host.	161

Table 6.23	<p> SNAK is used to reduce the amount of traffic to be analyzed. 163 For Normal host file the reduction was with a percentage of 60%. </p>
Table 6.24	<p> The Outbound scanning result of the monitored host. 165 </p>
Table 6.25	<p> Reduction rate resulted by layer 2. 165 </p>
Table 6.26	<p> LBDF Results for the cases 1 – 5. 167 </p>
Table 6.27	<p> The detection rates after applying Dataset15. 167 </p>
Table 6.28	<p> Results of methods evaluation. 168 </p>
Table B.1	<p> SNAK is used to reduce the amount of traffic to be analyzed. 193 For Normal host file the reduction was with a percentage of 60%. </p>
Table B.2	<p> The Outbound scanning result of the monitored host 196 </p>
Table B.3	<p> SNAK is used to reduce the amount of traffic to be analyzed. 197 For this file the reduction was with a percentage of 84%. </p>
Table B.4	<p> The list of candidate peaks resulted by the periodogram 199 </p>
Table B.5	<p> The verified peaks by ACF, only one peak is considered and 199 the rest are discarded. </p>
Table B.6	<p> The Outbound scanning result of the monitored host 200 </p>
Table B.7	<p> LBDF summarized results for the Dataset15 files. 201 </p>
Table B.8	<p> LBDF results with Dataset15 files 202 </p>
Table B.9	<p> SNORT results with Dataset15 files 203 </p>
Table B.10	<p> Rishi results with Dataset15 files. 204 </p>
Table B.11	<p> Bothunter results with Dataset15 files. 205 </p>
Table B.12	<p> Strayer results with Dataset15 files. 206 </p>
Table B.13	<p> Livadas results with Dataset15 files. 207 </p>

Table B.14	Al-sadhan results with Dataset15 files.	208
Table B.15	The summary results of the selected detection methods.	209

LIST OF FIGURES

Page

Figure 1.1	Global Bot Infection, Bots are a global problem. The map shows the geographic locations of active bots at October 2010.	2
Figure 1.2	Thesis Scope	7
Figure 1.3	Research Framework	8
Figure 2.1	Botnet Overview	11
Figure 2.2	A typical Centralized Botnet structure	13
Figure 2.3	A typical Peer-to-Peer Botnet Structure	14
Figure 2.4	A typical Unstructured Botnet structure	15
Figure 2.5	List of vulnerabilities among some of the well known softwares.	17
Figure 2.6	The growth of Web Malware for the years (2006-2010).	20
Figure 2.7	Total Mobile Operating System Vulnerabilities – 2006-2010, IBM.	20
Figure 2.8	Botnet Life cycle	23
Figure 2.9	A simplified figure representing a signature-based detection model.	24
Figure 2.10	Two periodic signals with frequency of 5 Hz.	35
Figure 2.11	Figure 2.1 Examples of finite duration non-periodic signals.	36
Figure 2.12	Example of infinite duration non-periodic Signals.	36
Figure 2.13	Left plots show non-periodic signal with a period of 0.2 second, and its periodogram. Right plots show periodic pulses with a period of 0.2 second, and its periodogram with	38

	a large peak at the fundamental frequency ($f=5$ Hz).	
Figure 2.14	The periodogram shows P1 more significant than P2 while the ACF shows P2 to be more significant than P1.	41
Figure 2.15	Periodogram and Autocorrelation function will be used to detect sequence periodicity and locate the fundamental peak(s).	43
Figure 2.16	Packet capturing process based on Libpcap library	44
Figure 2.17	General Network of Traffic Monitoring System.	45
Figure 3.1	Architecture overview of the proposed detection framework.	55
Figure 3.2	Layered approach of the proposed framework.	56
Figure 3.3	Packet capturing process.	57
Figure 3.4	The generation of the network traffic representative.	58
Figure 3.5	Captured traffic is decomposed into (SNAK) and (data).	59
Figure 3.6	SNAK is obtained from packets that contain SYN or ACK in its header.	59
Figure 3.7	The general model used to extract SNAK packets and to generate the discrete time sequences $x[n]$.	61
Figure 3.8	The effect of the chosen aggregation intervals on the SNAK-data relationship.	63
Figure 3.9	The Sliding time-window mechanism.	64
Figure 3.10	Effect of window-size on the correlation results.	65
Figure 3.11	The general Model for testing periodical behavior in $x[n]$.	67
Figure 3.12	The general algorithm for testing periodicities of $x[n]$.	68
Figure 3.13	Left plots show non-periodic signal and its one sided periodogram that consists of many peaks. Right plots show	69

periodic train of rectangular pulses and its one sided
periodogram.

Figure 3.14	The general algorithm for testing periodicities of $x[n]$.	70
Figure 3.15	The general Model for testing malicious behavior of $x[n]$.	71
Figure 3.16	TCP scanning model.	72
Figure 3.17	Inbound Scan, scanning activity is performed against internal hosts.	73
Figure 3.18	The calculation of scan-warning threshold S_{th} for one host.	74
Figure 3.19	The calculation of the average S_{th} .	75
Figure 3.20	Outbound Scan, when internal host try to perform the scanning activity.	76
Figure 3.21	The general algorithm that will be used to detect scanning behavior.	78
Figure 3.22	The analyzer rule is to correlate hosts that perform scan activity and their traffic exposes a periodic behavior.	79
Figure 4.1	An overview of the bot detection framework.	83
Figure 4.2	The model that was used to capture the traffic of the monitored network.	84
Figure 4.3	Raw packets are processed in layer 2 to produce SNAK and DATA packets.	85
Figure 4.4	Filtering rules used to generate SNAK traffic (SNAK Filtering Rules).	85
Figure 4.5	The algorithm that shows how the Discrete Time sequence $x[n]$ is generated.	87
Figure 4.6	The algorithm that will be used to show the similar traffic	89

	behavior between SNAK and data traffic sequences.	
Figure 4.7	moving from time domain to frequency domain is done by calculating PSD of $x[n]$, the result of this process is FFT $x[n]$, which will be known as $X[k]$.	90
Figure 4.8	Evaluating periodogram of $X[k]$.	92
Figure 4.9	The algorithm that will be used to specify the candidate peaks.	94
Figure 4.10	The algorithm that computes ACF of $x[n]$.	95
Figure 4.11	The location of the two segments on the ACF around the candidate peak, the linear regression is also shown as a black thick line.	97
Figure 4.12	The slope of $S1$ and $S2$ are used to detect the state of curvature around the candidate peak.	98
Figure 4.13	The segments slopes define whether the ranges of selected values are concaving up or down, i.e. to show that they are Valley or Hill.	98
Figure 4.14	The process of using segmentation and linear regression to locate the verified peaks.	99
Figure 4.15	The fundamental frequency is calculated after locating the maximum value (peak) of ACF in the selected range $[R1, R2]$.	100
Figure 4.16	The algorithm that will locate the fundamental frequency F .	101
Figure 4.17	Periodogram is been used to extract candidate periods. These candidates are then verified against the autocorrelation.	102
Figure 4.18	The algorithm used to calculate Whs and Wls .	104

Figure 4.19	The algorithm used to calculate Sth.	105
Figure 4.20	Inbound Scan Detection algorithm.	106
Figure 4.21	S1 represents the accumulative account of scanning activities that every host performs.	107
Figure 4.22	The algorithm that will be used to evaluate the outbound scan rate (S1).	108
Figure 4.23	The connection failure rate will be compared against selected threshold value so that to distinguish between normal and scanning hosts.	109
Figure 4.24	The format of the obtained data from layer three, for periodic field it will contain either the value (1) for periodic sequences or (0) for non periodic.	110
Figure 4.25	The format of information from layer Four.	111
Figure 4.26	A list of the monitored hosts is generated to show the status of each host.	111
Figure 5.1	The capturing tool is used to forward the raw captured traffic to LBDF.	114
Figure 5.2	Packet counts of bidirectional traffic of NAv6 dataset file 1.	117
Figure 5.3	Sinusoidal signal $x(t) = \sin(2\pi t)$ with frequency = 5 Hz.	117
Figure 5.4	Outbound Scan Performed on many ports in one machine.	121
Figure 5.5	Outbound Scan Performed on one port in many machines.	121
Figure 5.6	Outbound Scan Performed on many ports in many machines.	122
Figure 5.7	The initiated test bed that is built to test LBDF.	123
Figure 5.6	Summarization of experiments used to validate and evaluate LBDF.	129

Figure 6.1	A snapshot of the captured raw-traffic.	131
Figure 6.2	Packet counts of bidirectional traffic of NAv6 dataset file 1 (Top), file 2 (Bottom), attack free. Aggregation interval = 10 seconds.	132
Figure 6.3	Time variation of the cross-correlation function between SNAK and data, packet counts of NAv6 dataset file 1 (Top) and file 2 (Bottom).	134
Figure 6.4	(TOP) Sinusoidal signal $x(t) = \sin(2\pi t)$. (Bottom) the resulted periodogram with a peak appeared at $f = 5\text{Hz}$.	136
Figure 6.5	The resulted ACF of the signal $x(t) = \sin(2\pi t)$.	136
Figure 6.6	The two segments S1 and S2 are located based on the ranges R1 and R2.	138
Figure 6.7	The location of the maximum ACF specifies the location of the fundamental period that is 200 milliseconds.	140
Figure 6.8	Periodogram specifies candidate peaks, the ACF and the linear regression were used to distinguish between fake and real peaks.	141
Figure 6.9	The top Figure shows the plot of solar radiation values, while the bottom Figure shows the periodogram of solar radiation data.	142
Figure 6.10	The correlator is used to match the results of layers 3 and 4.	150
Figure 6.11	A snapshot of the captured raw-traffic of the Zeus Bot.	151
Figure 6.12	The discrete time sequences $X[N]$ of Zeus captured traffic. The chosen aggregation interval = 0.5 second.	152
Figure 6.13	The resulted periodogram of $X[N]$ of the Zeus Bot, with a	153

	large peak occurring at $f = 0.0175$ Hz.	
Figure 6.14	A snapshot of the captured raw-traffic of the IRC Bot.	155
Figure 6.15	The resulted periodogram of $X[N]$ of the IRC Bot, with a large peak occurring at $f = 0.11$ Hz.	156
Figure 6.16	A snapshot of the captured raw-traffic of the Conficker Bot.	157
Figure 6.17	The resulted periodogram of $X[N]$ of the Conficker Bot, with a large peak occurring at $f=0.285$ Hz.	158
Figure 6.18	A snapshot of the captured raw-traffic of a scanner host.	159
Figure 6.19	The discrete time sequences $X[N]$ of Scanner Host captured traffic. The chosen aggregation interval = 1 second.	160
Figure 6.20	The resulted periodogram of $X[N]$ of the Scanner Host that consist of many small peaks.	161
Figure 6.21	A snapshot of the captured raw-traffic of the Normal host.	162
Figure 6.22	The discrete time sequences $X[N]$ of Normal1 host captured traffic. The chosen aggregation interval = 1 second.	163
Figure 6.23	The resulted periodogram of $X[N]$ of the Normal host that consist of many peaks.	164
Figure 6.24	LBDF results of the selected files for Case 1 – 5 described in section 5.3.	166
Figure 6.25	Comparison between the evaluation results of each method.	168
Figure A.1	Packet counts of bidirectional traffic of NAv6 dataset file 1 (Top), file 2 (Bottom), attack free. Aggregation interval = 10 seconds.	182
Figure A.2	Packet count of bidirectional traffic for Monday week 1 (Top) and Tuesday week 1 (Bottom), of the DARPA 1999	182

	dataset, No attacks. Aggregation interval = 600 seconds.	
Figure A.3	Packet count of bidirectional traffic for Monday week 1 (Top) and Tuesday week 1 (Bottom), of the DARPA 1999 dataset, No attacks. Aggregation interval = 600 seconds.	184
Figure A.4	Packet counts of bidirectional traffic for Monday week 2 (Top) and Tuesday week 2 (Bottom) of the DARPA 1999 dataset, with attacks. Aggregation interval = 600 seconds.	185
Figure A.5	The effect of Clia1 attack on the relation between SNAK and data traffic sequences of NAv6 dataset file 4. Aggregation interval = 10 seconds.	186
Figure A.6	The effect of Flood1 attack on the relation between SNAK and data traffic sequences of NAv6 dataset file 6. Aggregation interval = 10 seconds.	187
Figure A.7	Time variation of the cross-correlation function between SNAK and data packet counts of DARPA 1999 dataset week 1, Monday (Top) and Tuesday (bottom).	188
Figure A.8	Time variation of the cross-correlation function between SNAK and data packet counts of NAv6 dataset file 3 with injected abnormal instances.	190
Figure A.9	Time variation of the cross-correlation function between SNAK and data, packet counts of DARPA week 2 dataset, Monday (Top) and Tuesday (bottom).	191
Figure A.10	Plots of Nav6 File 4 with clia1 attacks.	192
Figure B.1	A snapshot of the captured raw-traffic of the Normal 2 host.	193
Figure B.2	The discrete time sequences $X[N]$ of Normal 2 host captured	194

	traffic. The chosen aggregation interval = 1 second.	
Figure B.3	the resulted periodogram of $X[N]$ of the Normal host that consist of many small peaks.	195
Figure B.4	A snapshot of the printer server captured raw-traffic.	197
Figure B.5	The discrete time sequences $X[N]$ of printer server captured traffic. The chosen aggregation interval = 1 second.	198
Figure B.6	the resulted periodogram of $X[N]$ of the Printer Server, with many small peaks and a candidate peak occurring at $f = 0.3333$ Hz.	198

LIST OF ABBREVIATIONS

ACF	Auto Correlation Function
ACK	Acknowledgment
ASCII	American Standard Code for Information Interchange
bpp	Byte Per Packet
bps	Byte Per Second
C&C	Command and Control
DNS	Domain Name Service
DDNS	Dynamic Domain Name System
f	Frequency
FFT	Fast Fourier Transform
FIN	Finish
FN	False Negative
FP	False Positive
fph	Flow Per Hour
FTP	File Transfer Protocol
HTTP	Hyper Text Transfer Protocol
ICMP	Internet Control Message Protocol
ICQ	Instant Messaging Computer
IDS	Intrusion Detection System
IM	Instant Messaging
Ipsrc	The IP Source
IPv4	IP Address Version 4
IRC	Internet Relay Chat

ISP	Internet Service Provider
LBDS	Layered Approach Detection System
LAN	Local Area Network
MAN	Metropolitan Area Networks
MSN	Microsoft Networks
NXDOMAIN	Non Existent Domain
P2P	Peer to Peer
PCF	Partial Creation Filter
PSD	Power Spectral Density
RST	Reset
SMS	Short Message Service
SSL	Secure Socket Layer
SYN	Synchronize
TCP/IP	Transmission Control Protocol\Internet Protocol Stack
TN	True Negative
TP	True Positive
VoIP	Voice over Internet Protocol
WAN	Wide Area Network
Whs	Weight of the High Severity Ports
Wls	Weight of the Low Severity Ports
WWW	World Wide Web

LIST OF APPENDICES

Appendix A	LBDS Detection Results	174
Appendix B	LBDS Evaluation Results	191

RANGKA KERJA PENGESANAN BERLAPIS BOTNET BERDASARKAN PEMROSESAN ISYARAT DAN ANALISIS MASA DISKRET

ABSTRAK

Transaksi kewangan dalam talian dan maklumat sensitif yang banyak saling bertukar di Internet. Ini mengalih tumpuan penyerang siber daripada perasaan ingin tahu kepada mendapatkan keuntungan kewangan. Penyerang menggunakan perisian hasad yang berbeza untuk mencapai matlamat mereka. Botnet dianggap antara perisian hasad yang berbahaya kerana kuasanya yang mampu mengawal pelbagai mesin dan memberi ancaman kepada pengguna Internet.

Tesis ini membentangkan suatu pendekatan baru dalam bidang pengesanan botnet. Ia memperkenalkan rangka kerja baru yang dinamakan Rangka Kerja Pengesanan Botnet Berlapis (LBDF) yang dapat mengesan rakan botnet dengan berkesan. Rangka kerja ini berfungsi dalam domain kekerapan dan bukannya dalam masa domain. LBDF dilengkapi dengan algoritma pengesanan 'pengimbasan-hasad'. Algoritma LBDF menggunakan peraturan SYN, ACK (SNAK) untuk mengurangkan jumlah kesesakan rangkaian dan menukarkan trafik yang telah dikurangkan menjadi sampel data diskret. Kemudian LBDF mengaplikasikan kedua-dua periodogram dan fungsi autokorelasi membulat bagi mengesan sebarang keberkalaan tersembunyi di dalam jujukan sampel. Jika pelakuan berkala dikesan, kekerapan jujukan dan alamat IP komputer akan direkodkan. Oleh itu, alamat komputer peribadi dengan pelakuan berkala akan disimpan ke dalam pangkalan data dan dilabelkan sebagai mencurigakan. Jika mana-mana mesin yang mencurigakan menunjukkan pelakuan

pengimbasan-perosak, ia akan diisytiharkan sebagai bot. Bot yang mempunyai ciri yang sama dikumpulkan sebagai jenis botnet yang sama. Walaupun LBDF tidak mengesan bot yang tidak berkala atau tidak aktif contohnya, tidak berkomunikasi dengan master bot atau melakukan sebarang tindakan, ia akan mengesan mereka apabila mereka menunjukkan pelakuan yang mencurigakan. Pendekatan ini adalah berbeza dengan pendekatan lain kerana ia tidak terhad kepada protokol yang spesifik kepada protocol C&C (contohnya; HTTP, IRC) atau struktur botnet yang spesifik (contohnya; P2P, Berpusat) atau pelakuan serangan (iaitu; SPAM, DDOS) yang tidak memerlukan sebarang pengetahuan terdahulu bot yang dikesan.

Penilaian LBDF menunjukkan bahawa algoritma pengesanan adalah tepat, pantas dan berskala jika dibandingkan dengan rangka kerja pengesanan yang ada. LBDF mampu mengesan P2P, HTTP, IRC, bot berpusat atau yang tidak berstruktur. Justeru, prestasi LBDF F-measure adalah 26% lebih baik berbanding rangka kerja pengesanan botnet yang lain. Hasil daripada pengaplikasian algoritma pengurangan trafik rangkaian yang diadaptasikan oleh LBDF, kadar pengurangan dalam trafik yang dikaji adalah dalam julat 20%-90%, pengurangan ini meningkatkan prestasi LBDF dan meningkatkan daya pemprosesan tanpa menjejaskan matlamat utama LBDF.

LAYERED BOTNET DETECTION FRAMEWORK BASED ON SIGNAL PROCESSING AND DISCRETE TIME ANALYSIS

ABSTRACT

A massive volume of online financial transactions and sensitive information is exchanged over the Internet. This has shifted the focus of cyber attackers from curiosity to financial gain. Attackers use different malware to achieve their goals. Among the various forms of malware; the botnet is considered as the worst, because of its vast computing power, ability to control many machines and its significant threat to the Internet users.

This thesis presents a new approach in the area of botnet detection. It introduces a framework called Layered Botnet Detection Framework (LBDF) that can detect botnet members efficiently. This framework works in the frequency domain rather than in the time domain. LBDF is equipped with a ‘malicious-scanning’ detection algorithm. The LBDF algorithm uses SYN, ACK (SNAK) rules to reduce the volume of network captured traffic and to convert the reduced traffic into discrete time sequences. Then LBDF applies both a periodogram and circular autocorrelation function to these sequences to detect any hidden periodicities. If periodic behavior were detected, the frequency of the sequence and the IP address of the monitored computer will be recorded. Thus the IP address of PCs with periodic behavior will be saved in a database and labeled as suspicious. If any of the suspicious machines performs a malicious-scanning action, it will be declared as a bot. Bots that have similar features are grouped together as members of the same botnet.

Although LBDF does not detect bots that are non-periodic or inactive i.e. not communicating with their bot master or performing any action, it will detect them as soon as they exhibit suspected bot behavior. The proposed approach is different than other approaches, since it is not limited to specific C&C protocols (e.g., HTTP, IRC) or to specific botnet structures (e.g., P2P, Centralized) or attack behaviors (i.e. SPAM, DDOS), neither does it require any prior knowledge of the detected bots.

The evaluation of LBDF shows that the detection algorithm is accurate, fast and scalable compared to existing bot detection frameworks. LBDF is capable of detecting P2P, HTTP, IRC, centralized and even unstructured bots. In this respect, the LBDF F-measure is better by 26% compared with other botnet detection frameworks. As a result of applying the network traffic reduction algorithm adopted by LBDF, the reduction rate in the analyzed traffic was in a range of 20% - 90%, this reduction improves the performance of LBDF and increases its throughput without affecting the main goal of LBDF.

CHAPTER ONE

INTRODUCTION

A few years ago, protecting a computer system or networks was mainly required to prevent the threat from viruses and worms. Nowadays, the situation has changed drastically; the biggest threat faced by network hosts is malware, which is written by cyber programmers with the intent of malicious activities. Malware may take the form of viruses, worms, Trojans, botnet or other malicious programs.

Among the various forms of malware, botnet is considered as the most serious means for conducting online crimes (FBI, 2011b). This threat is triggered due to its large scale and geographical diversity of the network hosts enlisted in a Botnet. The large number of enlisted bots; gave the Botnet its vast computing power (Guofei Gu, Roberto Perdisci, Junjie Zhang, & Wenke Lee, 2008). This vast computing power coupled with the easy controlling of botnet from anywhere in the world; makes the botnet a powerful cyber weapon and an effective tool for performing malicious activities. Botnets become sophisticated more and more every day by employing variety of techniques (e.g., sophisticated executable packers, rootkits, protocol evasion techniques, such as moving away from IRC and taking control of, HTTP, VoIP, IPV6, ICMP, Skype protocols, etc). Bots are more evasive to signature based detection systems, anomaly-based detection systems as well as DNS and data mining based intrusion detection systems. These evasion techniques; improve the survivability of botnets and the success rate of compromising new hosts. Additionally, botnets have also added (and continue to add) new mechanisms to hide traces of their communications.

1.1 The scale of Botnet problem

A bot is defined as a computer that is compromised by malicious software that enables a remote computer to control it. Bots are part of a network of infected machines, known as a “botnet” that spread globally as shown in Figure 1.1. The process of estimating botnet size and calculating the botnets population is a tedious task; the size and growth of botnets differ widely. For example, the Mariposa botnet (FBI, 2011a), it contains 12.7 million infected computers, while Zeus (Binsalleeh, et al., 2010) has more than 1,400 command and control servers with undetermined numbers of infected hosts. The figures are scary. The total number of zombies is near to 60 million. Table 1.1 shows some evaluations of the number of active bots at the end of 2010 according to Message Labs Intelligence 2010 Annual Security Report (Symantec, 2011).

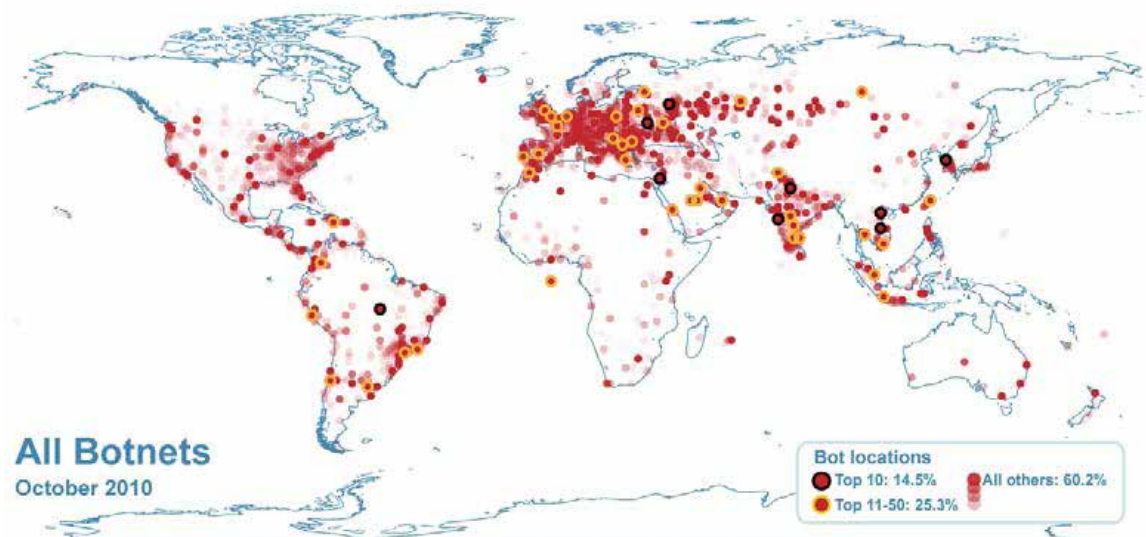


Figure 1.1 Global Bot Infection, Bots are a global problem. The map shows the geographic locations of active bots at October 2010.(Symantec, 2011)

Table 1.1 State of Botnet at the end Of 2010 (CISCO, 2011)

Botnet	Est.Botnet size	Country of infection
Rustock	1100k to 1700k	USA (17%), Brazil (7%), India (7%)
Grum	310k to 470 k	Russia (12%), India (8%), Vietnam (8%)
Cutwail	560k to 840k	India (17%), Russia (16%), Ukraine (8%)
Maazben	510k to 770k	Russia (11%), India (10%), Brazil (7%)
Mega-D	80k to 120k	Russia (15%), Ukraine (14%), Brazil (7%)
Cimbot	32k to 48k	Italy (27%), Spain (25%), France (14%)
Bobax	250k to 370k	India (32%), Russia (25%), Ukraine (9%)
Xarvester	17k to 25k	Italy (15%), UK (10%), Poland (8%)
Festi	8k to 12 k	Vietnam(24%), Indonesia(21%), India
Gheg	8k to 12 k	Spain (12%), Indonesia (21%), India (10%)
Unnamed	490k to 740k	
other	220k to 340k	
Total	3500k to 5400k	India (9%), Russia (9%), USA (7%)

1.2 Research Motivation

Regardless of how malware reaches a computer, the challenge is to identify the infected machine and heals it as soon as possible before any harm is caused. The past recent years are witnessed of different approaches that have been proposed to detect botnets and to combat their threat against cyber-security, but these approaches were based on a specific part of botnet lifecycle like, scan, spam, etc. or a specific abnormal behavior of a network traffic or, a specific communication protocols like, IRC, P2P and HTTP that are used by botnet Command and Control servers (C&C) or a certain topology e.g, centralized. All of these properties are specific properties and it is not necessary that all types of botnets contain it. Therefore, previous methods are suitable only for specific botnet type or structure. Diversity of botnet protocols and different structures; make botnet detection a very challenging task.

1.3 Problem Statement

One of the most critical issues in Cyber Security is the botnet detection problem. Bots are stealthy in nature and usually do not aggressively consume CPU/memory/bandwidth resources, or perform noticeable damage to computers, such as disabling existing antivirus. Thus, a host-based solution method that is very specific to a certain botnet's structure or a certain communication protocol is not desirable because:

- Bots are flexible in their nature.
- Continuously evolving with flexible design.
- Different protocols and structures are used to organize and control the botnet.
- Bot life cycle consists of several different stages and aspects that developed and changed continuously.

That is why many existing Bot detection techniques become ineffective, as bots change their structure or C&C techniques. Despite the concerted efforts given in the literature, diversity of botnets protocols and structures makes botnet detection a challenging task and unsolved problem for the online community (ENISA, 2011; FBI, 2011a; IBM, 2010). Botnet detection problem can be solved through the detection of the command and C&C communication channels and the host's malicious-activities that is proposed in this thesis.

1.4 Research Objective

The objective of this thesis is to construct a new approach to detect botnet members in the monitored network. This approach is independent of botnet C&C protocols and structures. In addition, it does not require any priori knowledge (signature) of bots. It is assumed that the detection of the periodic C&C communication channels traffic together with the detection of the malicious-scan activities makes it possible to detect botnet members in the monitored network.

Therefore, the main objectives of this thesis are:

- To create a traffic representative that functions in frequency domain.
- To detect bots, independent of the bots structure and the communication protocols used.
- To evaluate the performance and accuracy of the proposed framework compared to other existing frameworks.

1.5 Thesis Contribution

Computer users, applications and bots utilize network in the same manner, but with different intentions. The naïve nature of users and applications activities differs than the malicious activities performed by bots. The proposed framework should be able to distinguish the normal traffic caused by a legitimate user or applications from the malicious traffic caused by bot activity.

The main expected contribution of this thesis is to propose and to design a new bot detection model. However, this thesis contribution summarizes as follow:

- Ø **Traffic reduction**, to introduce a new technique that will be able to create a true representative of the monitored network traffic with a discrete time sequences.
- Ø **An algorithm that computes the PSD of the resultant Discrete Time Sequences**, an algorithm that can be used to understand, and to model the normality of the network traffic in the frequency domain rather than the time domain.
- Ø **A bot detection model with low false positives**, a new detection model that is capable of identifying the existence of all known and unknown types of bots, independent of both, the botnet structure and the used communication protocols.
- Ø Enhanced model in terms of **performance and accuracy**, compared to the existing models.

1.6 Thesis Scope

The scope of this work as shown in Figure 1.2 is limited to Inbound/Outbound, IPv4 and TCP traffic, captured from the observed hosts in the monitored local area network (LAN). In the captured traffic, only details within the packet headers are of interest.

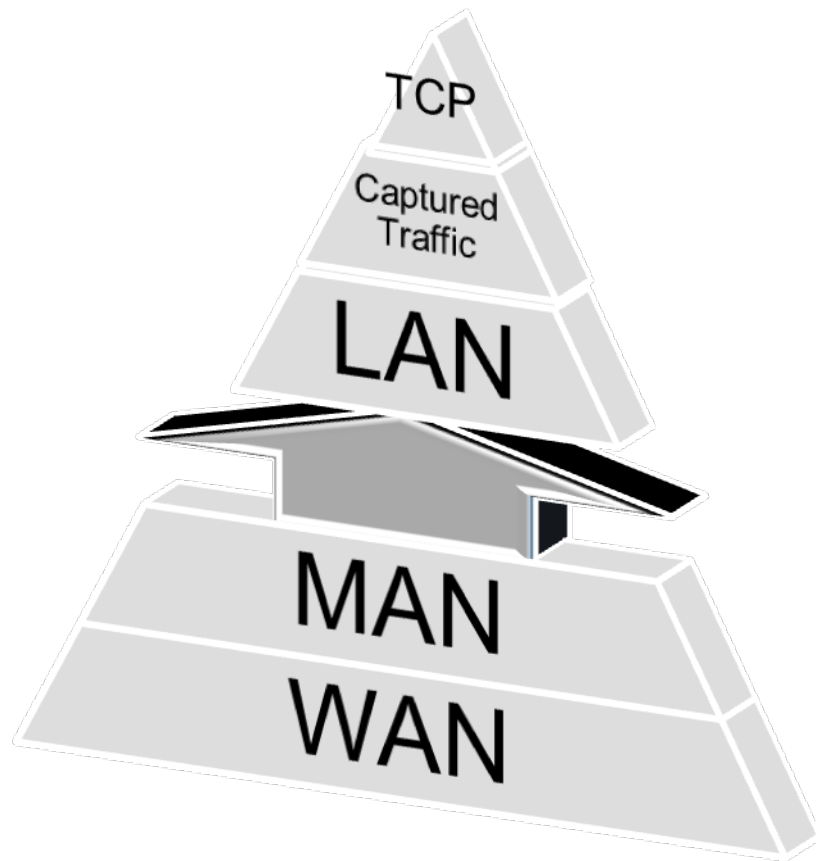


Figure 1.2 Thesis Scope

1.7 Research Framework

Figure 1.3 describes the complete research framework of this thesis.

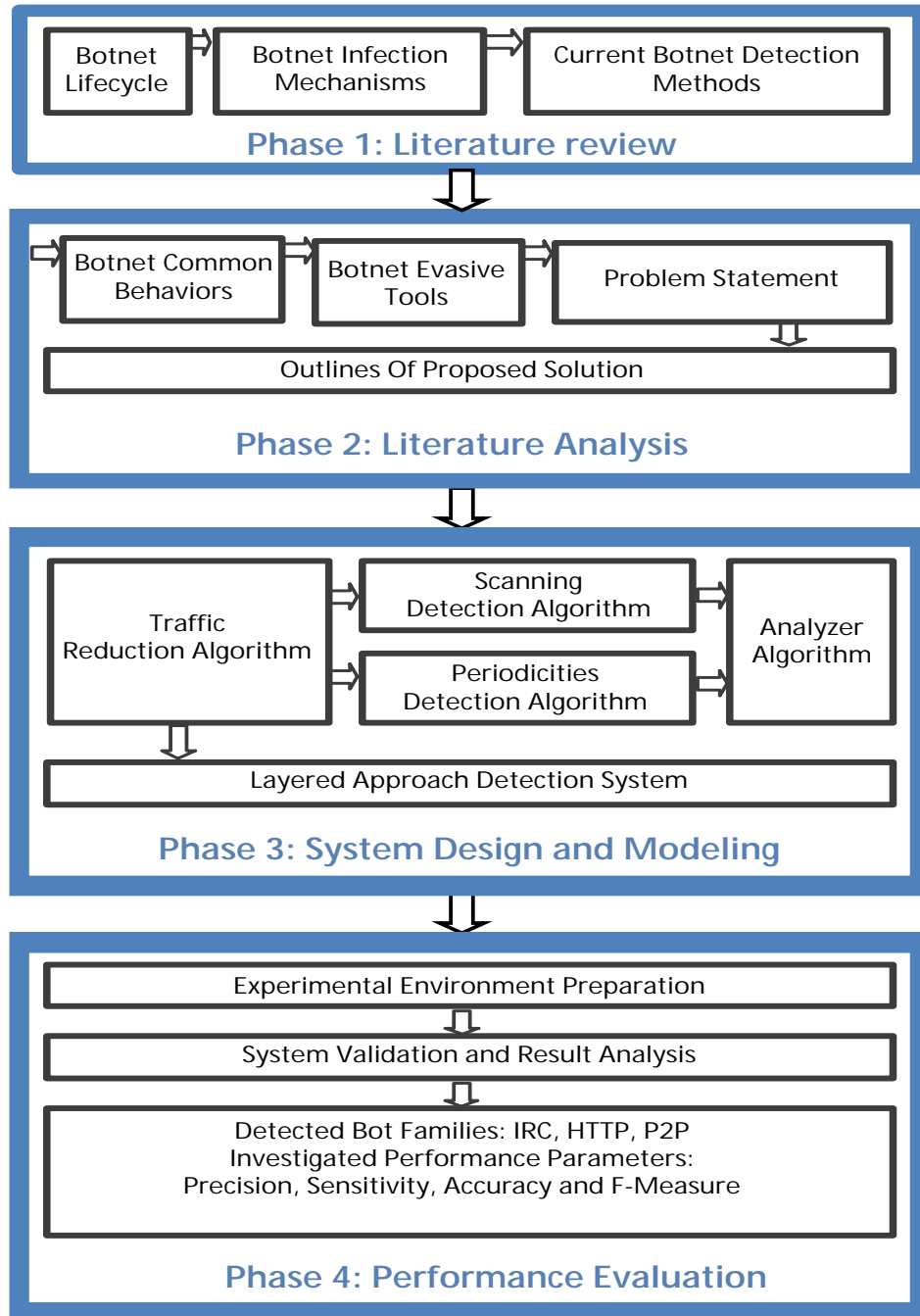


Figure 1.3 Research Framework

1.8 Thesis Outline

This thesis is organized into seven chapters. This chapter (**Chapter 1**) presents the objectives of this thesis. It starts by presenting a background discussion for the Bot problem along with our research objectives and contributions.

In **Chapter 2**, literature review will be presented along with some fundamental concepts related to this work and issues surrounding it. Other botnet detection models will be discussed, as well as the most current and related works related to botnet detection. This chapter also provided motivation for our work by describing some candidate architectures and the limitations of those proposed solutions.

Chapter 3 covers the methodology discussion on how the proposed solution was designed. The algorithm to shift captured traffic from time domain to frequency domain, and to accurately -detect and specify- the frequency (time) of the monitored signal (traffic) will be introduced in this chapter.

The implementation details and issues regarding the illustration of the detection model implementation were presented in **Chapter 4**. While the explanation of the performed experiments and the used datasets are presented in **Chapter 5**.

The results obtained by the experiments in Chapter 5 are the primary content of **Chapter 6**. Finally, **Chapter 7** presents three main headings conclusion, recommendation and the possible future work for this study.

CHAPTER TWO

LITERATURE REVIEW

Botnet threats are a magnified extension of previous computer threats, combined with C&C systems and capable of infecting hundreds of thousands of computer systems. Despite botnet recency, this area witnessed a significant number of researches and proposed solutions. In this chapter, an overview about botnets is introduced. Also, this chapter presents the botnet and the related researches including infection mechanisms, botnet communication protocols, C&C models, malicious behavior, previous and current, bot detection methods and botnet defense. Moreover, some related topics like periodic signals, periodograms and circular autocorrelation function, also will be discussed.

2.1 Introduction

The related work that will be discussed in this chapter concerns previous approaches of botnet detection methods; their pros and cons. In particular, Signature Based, Anomaly Based, DNS Based and Data Mining Based Techniques will be discussed. Primarily focus will be concentrated on the efforts that have been made to detect botnets, based on known, anomalous, preprogrammed, repetitive and correlated behavior of botnet members. These effort's advantages and disadvantage will be discussed individually. Moreover, a comparison between those models will be performed and some examples of each of them will be provided.

2.2 Background of Botnet

For better understanding of botnet, some key terms are introduced that are related to the botnet community. The most related topics to botnet detection are shown in the general outline presented in Figure 2.1.

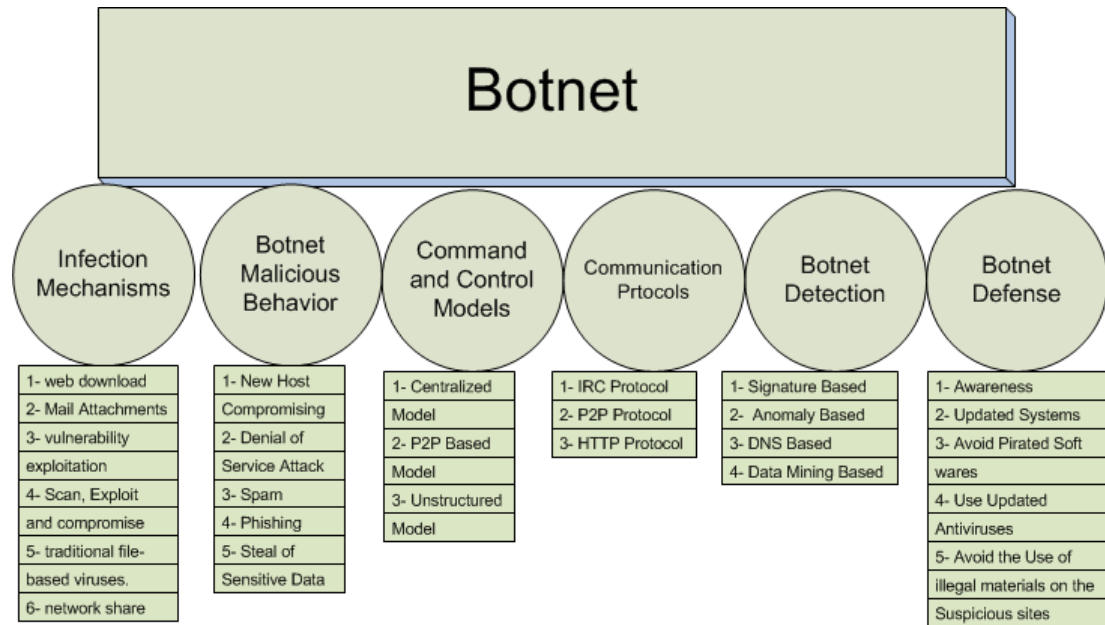


Figure 2.1 Botnet Overview

2.2.1 Definitions Related to Botnet

∅ **Bot** - the term bot comes from the word robot, which means "worker". In the world of computers, bot is a generic term, used to describe an automated process (Geer, 2005; Ianelli & Hackworth, 2005; Saha & Gairola, 2005). A bot is usually referred to as automated software, which is capable of performing certain predefined tasks repeatedly.

∅ **Botnet** - are a group of compromised computers (or zombies) that are under the control of a single entity called botmaster (Barford & Yegneswaran, 2007; Gu, Porras, Yegneswaran, Fong, & Lee, 2007; Saha & Gairola, 2005).

Ø *Command and Control (C&C)* is the commander channel that receives commands from the botmaster and conveys these commands to their bots in order to carry out various distributed and coordinated attacks remotely (B. AsSadhan, Moura, Lapsley, Jones, & Strayer, 2009; Bailey, Cooke, Jahanian, Xu, & Karir, 2009; Gu, Zhang, & Lee, 2008).

Bot infected computers can be controlled as:

- Directly, by initiating a connection with the infected computer known as the channel. Then, controlling it by commands hardcoded into the bot program, e.g. IRC and HTTP botnets.
- Indirectly, the bot initiates the connection with the control center/peer, sends a request and then performs the returned command e.g. P2P botnets.

2.2.2 Current and Expected Future Structure of Botnet

Botnet can be classified according to:

- Topology
- Communication protocols

✓ Botnet Classification According to Topologies

Botnets can be classified based on their C&C architectures as follows: (Chao, Wei, & Xin, 2009).

- a- *Centralized architecture*: in a centralized architecture, all bots are connected to a certain centralized C&C server, such as IRC and HTTP based botnets as shown in Figure 2.2. This architecture is considered as the easiest to construct

and implement; that is why this structure is still in use in the cyber world till now. However, this architecture suffers from the one single point of failure architectures, where it is easy to identify the C&C server and thus, bring down the whole botnet.

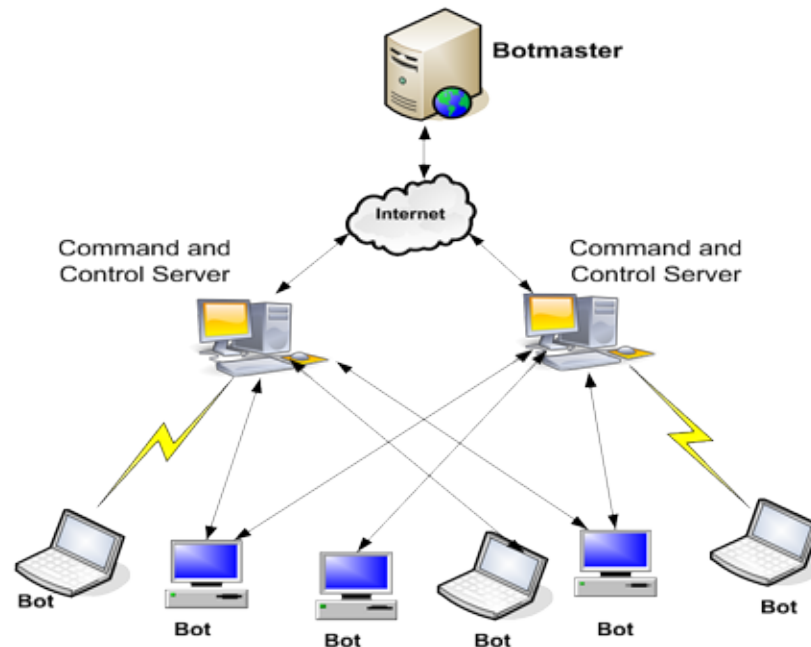


Figure 2.2 A typical Centralized Botnet structure.

- b- ***P2P or Decentralized architecture***: shown in Figure 2.3, in this architecture there is no centralized point for C&C. So that any node in the network can act as a client and as a server, P2P architecture employs the P2P protocols to present a various distributed C&C servers. This architecture is considered difficult to discover and destroy, due to the anonymity and the distributed nature of the P2P architecture (Grizzard, Sharma, Nunnery, Kang, & Dagon, 2007).

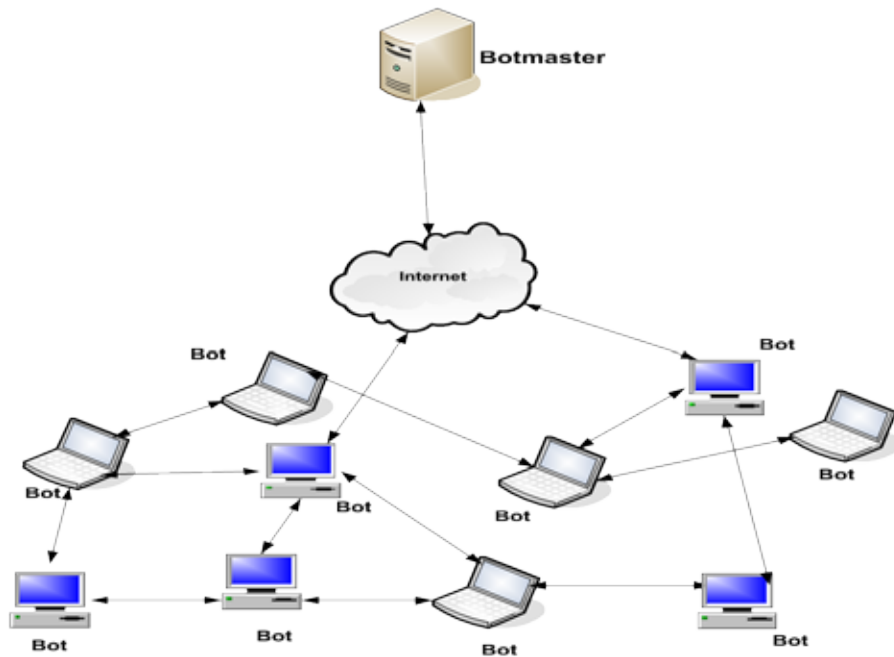


Figure 2.3 A typical Peer-to-Peer Botnet Structure

c- *Unstructured C&C architecture*: shown in Figure 2.4 and considered to be the extreme case of P2P botnets; where each bot is connected to one peer and doesn't know anything about other peers in the botnet, and more importantly, the bots in this structure are randomly organized (Clarke, Sandberg, Wiley, & Hong, 2001; Gnutella, March 2001). In this model there is no direct connection between the bot and the bot master; the bot master has to search the Internet and posts the required tasks to the bot when it finds one. Such a system is simple to design and to implement. The single bot detection would never compromise the whole botnet. However, this structure will not be effective as other structures; as it doesn't have a guarantee of delivery and suffer from high message latencies.

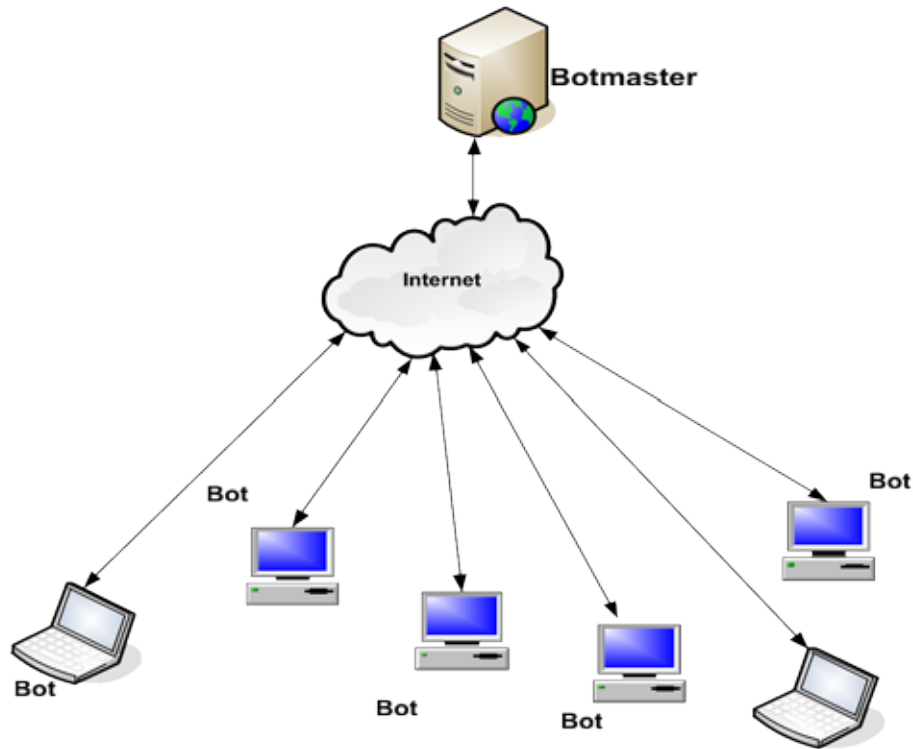


Figure 2.4 A typical Unstructured Botnet structure

The general properties of the different structures are summarized in Table 2.1

Table 2.1 C&C topologies and basic properties (Bailey, Cooke, Jahanian, Xu, et al., 2009)

Topology	Design Complexity	Detectability	Message Latency	Survivability
Centralized	Low	Medium	Low	Low
Peer-to-Peer	Medium	Low	Medium	Medium
Unstructured	Low	High	High	High

✓ Botnet Classification According to Communication Protocols

It is essential to have a communication channel between the bots and their owner, so that; the botnet owner can control his bots and send them the required commands.

Establishing these channels (connections) and maintaining them are based on network communication protocols.

Therefore, based on the used network protocols; botnets can be classified as (Tyagi & Aghila, 2011):

- a- **IRC-based:** in this botnet, bots are controlled via Internet Relay Chat (IRC) channels. IRC was the first and the common protocol used by botmasters, to send commands to the infected machines. IRC can be easily detected and network security devices can be configured to block IRC traffic.
- b- **IM-based:** which uses communication channels provided by instant messaging (IM) services such as AOL, MSN, and ICQ etc. IM is Low popular botnet communication channel; because it is difficult to create an individual IM accounts for each bot. Bots should be online all the time and keep connected to the network. IM services do not permit the same account, to log on to the system, from more than one host at the same time; each bot needs its own IM account; because automatic account registration is prevented in most of IM services. This will limit the number of registered IM accounts i.e. limits the number of concurrent online bots.
- c- **Web-based:** This technique is based on the popular communication HTTP protocol, which is difficult to be detected and can easily bypass network security devices. The botmaster controls his zombies from anywhere in the world through the World Wide Web by using the HTTP. The bot master defines a web server, the bot connects to the defined web server, receives commands and responses back to the same web server.
- d- **Other:** Botnets that use their own protocols to communicate, protocols that are based on the TCP/IP stack, i.e., that use transport-layer protocols such as TCP, ICMP and UDP.

2.2.3 How Bots Work

Bots spread across the Internet by exploiting vulnerabilities on unprotected computers to infect them and to report back to their master. Then the bots stay hidden until they are instructed to carry out another task. Based on 2010's results, the top 20 most common vulnerabilities were found in software developed by four companies: Microsoft, Adobe, Oracle and ACDSee, as shown in Figure 2.5.

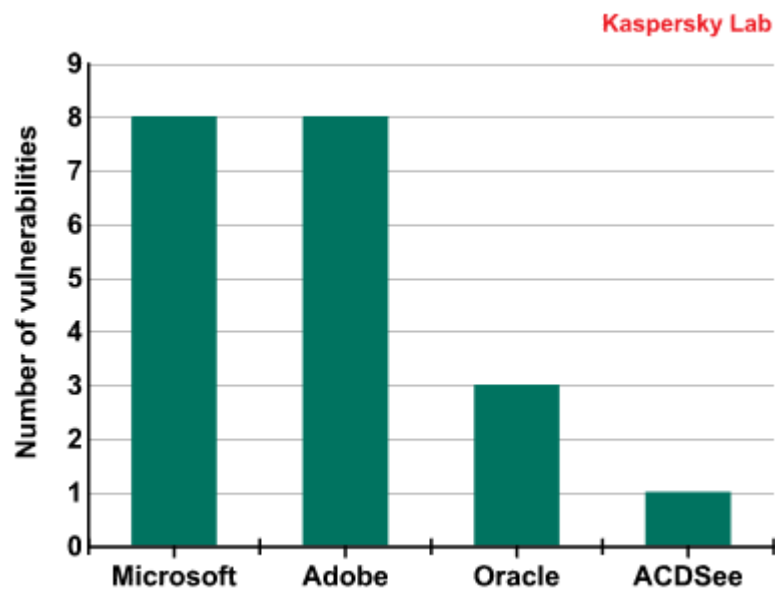


Figure 2.5 list of vulnerabilities among some of the well known softwares.

The compromised computer can be used to carry out a variety of automated tasks. Such as, Sending (Spam, Viruses, and Spyware), stealing personal information (credit card numbers, bank credentials, email address lists and Other sensitive personal information), launching denial of service (DoS) attacks, clicking on internet ads to boost up web advertising billings and extortion in which attackers ask to be paid, or they will attack the online services or the website of a certain company.

2.2.4 Botnet evolution

Like many services on the Internet, bots started as a useful tool without any malicious intent. Bots originally, were developed to sit on an IRC channel and perform several tasks, for its owner. Bots evolved from playing games with IRC users (GM: IRC bot, 1989, by Greg Lindahl), to a password stealer and backdoor (PrettyPark, 1999 'Trojan.PSW.CHV'), then it has the ability to remote control IRC clients by utilizing IRC vulnerabilities (SubSeven Trojan/Bot:By the late 1990s) (Tyagi & Aghila, 2011).

In 2000, Global Threat (GTBot) appeared, this bot can execute commands in response to events on the IRC server, and it supports raw TCP and UDP socket connections. GTbot had the capabilities of port scanning, flooding and cloning etc. (Jing Liu, Yang Xiao, Kaveh Ghaboosi, Hongmei Deng, & Jingyuan Zhang, 2009). In 2002, SDBot appeared which represents a new era in the evolutionary chain for bots with available source code, which made it accessible to many hackers. Moreover these types were easy to modify and to maintain (Jing Liu, et al., 2009). In 2002 bots with modular design appeared (Agobot, aka Gaobot, 2002), the modular design; allows the botmaster to update modules as new techniques or sites are available (Kola, 2008).

The available bot evolution techniques, lead to the creation of botnet that depends on unique characteristics, rather than depending on the original code, like Spybot and MyTob, for example (Polybot, March of 2004) has the capability to appear in many different forms. The use of hybrid, social engineering and spoofed e-mail addresses appeared with (Mytob, 2005). Botnets moved away from the original IRC Command

& Control channel, and began to communicate over HTTP, ICMP and SSL ports, Fast Flux network based on DNS servers and of course the P2P protocols. (Sinit, 2003) is an example of the early malicious Peer-to-Peer bot. In January 2007, the Trojan.Peacomm bot appeared; it was the most recently known peer-to-peer bot (Chao, et al., 2009).

After exploiting all available protocols, botnet developers turned their attention to network architecture, moving their botnets structure from the legacy classic architecture (i.e. a centralized structure with one or more C&C), into the dynamic P2P structure, which has no C&C, large botnet based on P2P architecture appeared in 2007. P2P botnet have attracted Bonet developer as well as botnet researcher.

2009 was characterized by the increased sophistication and the complex malicious programs that have rootkit functionality, for example the year 2009 witnessed (global epidemics, web attacks, web botnets, SMS fraud, the use of new platforms such as Mac OS and mobile operating systems and attacks on social networks) (Szor & Kaspersky).

2010, The Year of the Vulnerability (Bail, 2011), Web malware grew by 139 percent in 2010 compared to 2009 as shown in Figure 2.6. Numbers of used botnet-technologies have progressed dramatically. Botnet is constantly growing more and more complex, (e.g. Mariposa, Zeus, Bredolab, TDSS, Koobface, Sinowal and Black Energy 2.0 botnets); which are considered among the most sophisticated malware ever created. 2010 also witnessed compound efforts of law enforcement agencies, antivirus vendors and telecom providers in trapping of cybercriminal and illegal

services or business on the Internet. The P2P share was present to declare an increase in peer-to-peer (P2P) activity and again to focus more on the rule of P2P in the future direction of botnet.

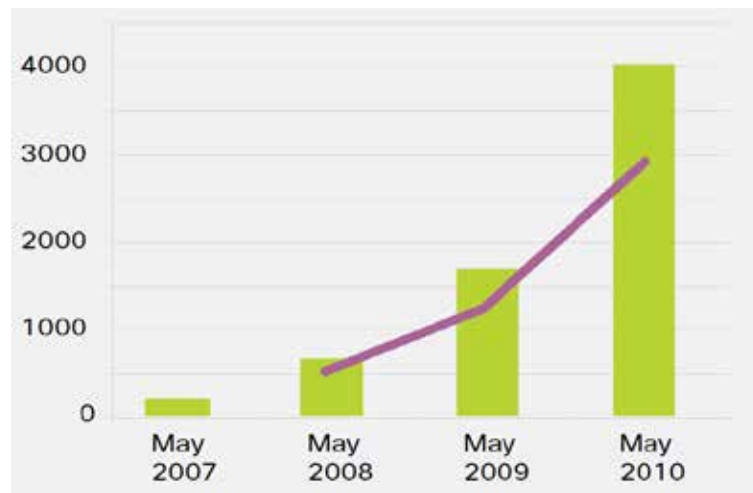


Figure 2.6 The growth of Web Malware for the years (2006-2010).

Also it was noticed that in the year 2010, attackers have shifted from internet and user's pc toward mobiles. A significant increase in mobile malwares gives us a black picture for the future of mobile botnets as shown in Figure 2.7.

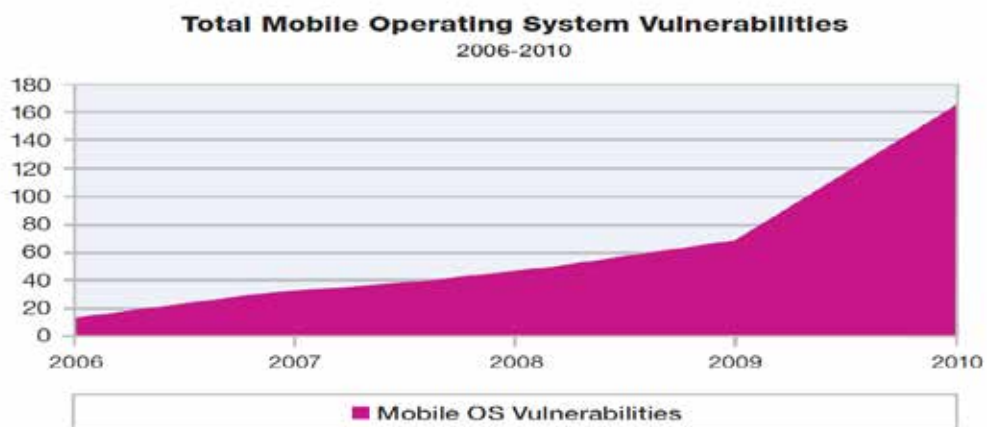


Figure 2.7 Total Mobile Operating System Vulnerabilities for the years 2006 to 2010, IBM

2.2.5 Botnet Potential Benefits

Botnets constructing have one or more of the following information processing, information harvesting and information dispersion. Information processing is used to process data, such as cracking passwords. While the use of information harvesting may includes obtaining financial data, password data, identity data, relationship data (i.e., email addresses, list of friends) or any available data on the host. Information dispersion includes providing false information from illegally controlled sources, creating denial of service attacks and sending out spam.

2.2.6 Botnet Infection Mechanisms

Botnets utilize many different infection mechanisms, such as employing malware (i.e. worms, trojan insertion), web driven-by download, mobile media, vulnerability exploitation, mail attachments, automatically scan-exploit-and-compromise, traditional file-based viruses, network share, as well as social engineering techniques and P2P file sharing networks, etc (Chao, et al., 2009).

2.3 Botnet Life Cycle & Detection Systems

Creating Botnet begins by sending a malware to vulnerable machines. Once vulnerability is found, the machine will be compromised; leading to the malicious bot binaries to be downloaded into the compromised hosts, turning it into a zombie (bot). This new bot in return will be redirected to a dynamic/static server address that is known for both the bot and his master. This server is known as a C&C server, where the botmaster (attacker) can login and issue commands to his bots to start an attack, scanning, infection...etc, (Chao, Wei et al. 2009). The most general phases in

Botnet lifecycle are: spread, infection, C&C, and attack, as shown in Figure 2.8.

Botnet life cycle events include:

- 1- Victim browses a website or clicks a link on email (e.g. phishing, drive-by download), then the browser is redirected to a malicious dropper site.
- 2- Victim is directed into downloading the dropper - or dropper is automatically downloaded through an exploit.
- 3- Dropper unpacks on the infected machine and runs.
- 4- Dropper informing its botmaster that it joined to the botnet.
- 5- The C&C secure the new client, sends encrypted malware with new instructions.
- 6- Dropper decrypted the malware and installs it. The dropper has to vanish by hiding, or delete itself so that users believe that no infection has occurred. Infected machine is turned into zombie (bot)
- 7- Malware contacts C&C, sends passwords/data/etc. as encrypted payload.
- 8- C&C updates the bot status and sends new instructions.
- 9- Bot responds by executing the commands and performing the required actions.
- 10- The bot contacts C&C sending its report.
- 11- In certain situations the bot is recommended to erase all commands and vanishes to remove any evidence on the botnet existence.

Steps 7, 8, 9 and 10 repeat indefinitely with the malware 'evidence' and C&C connection instructions changing constantly. The malware can be told to lay silent for a specified period of time.

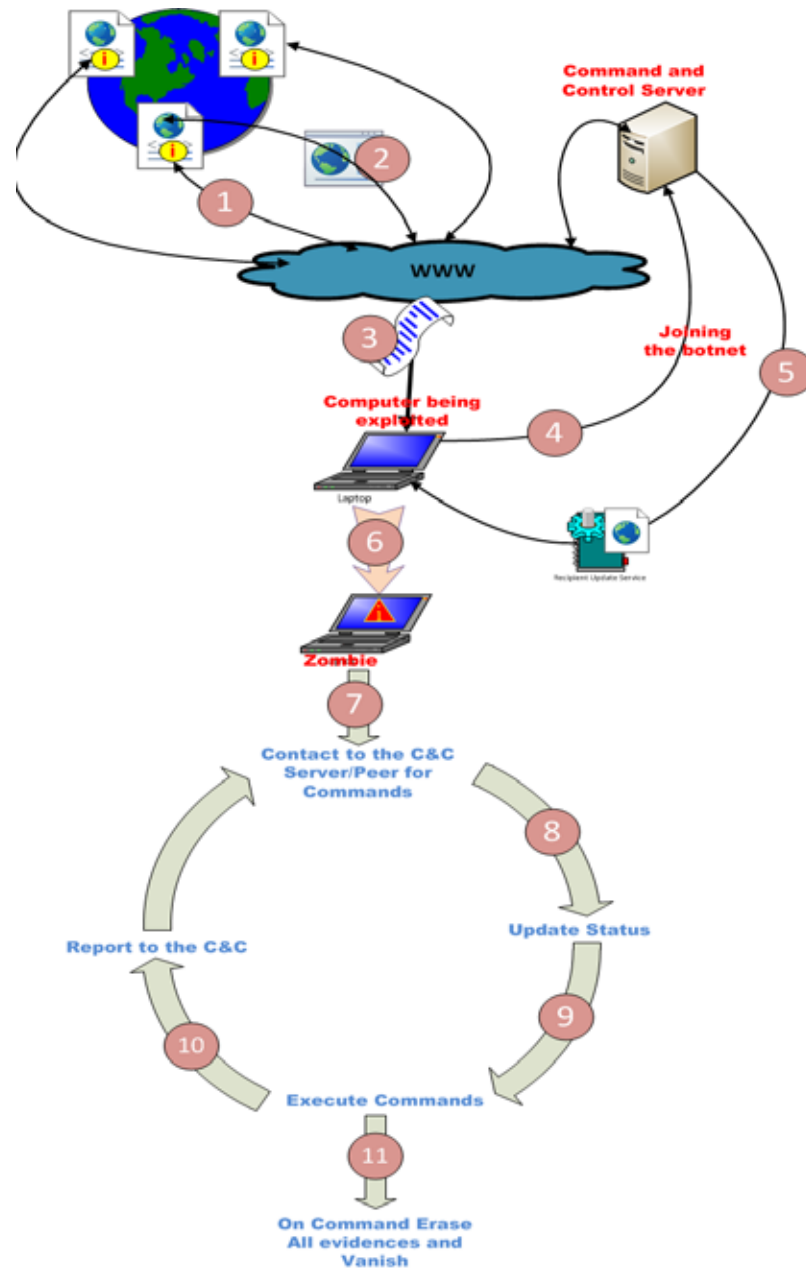


Figure 2.8 botnet life cycle.

Analyzing the malware after it enters the organization to obtain C&C details, can help removing it. Unfortunately, analyzing alone will not be enough as the infection lifecycle changes so quickly in a way that the analyzed malware no longer exists on the victim's machine.

2.4 Botnet Detection Methods

How the host is infected is not important as how to heal it, once the bot is created within the network hosts, the first priority is to identify the infected host and to heal it. Therefore, many efforts were initialized to detect botnets. The past recent years witnessed different approaches that have been proposed to detect botnets and to combat their threat against cyber-security. These approaches can be grouped into Signature Based, Anomaly Based, DNS based and Data Mining Techniques.

2.4.1 Signature-based Detection

Signature-based Detection explained in Figure 2.9 examines the network traffic for known patterns of a malicious activity; new types of attacks are not detected. Signature-based detection involves searching among network traffic, for a series of bytes or packet sequences or a set of attributes and matches them against a set of predetermined attribute lists. In case some particular network traffic has a match, the system has to alert administrators or to take a pre-defined action.

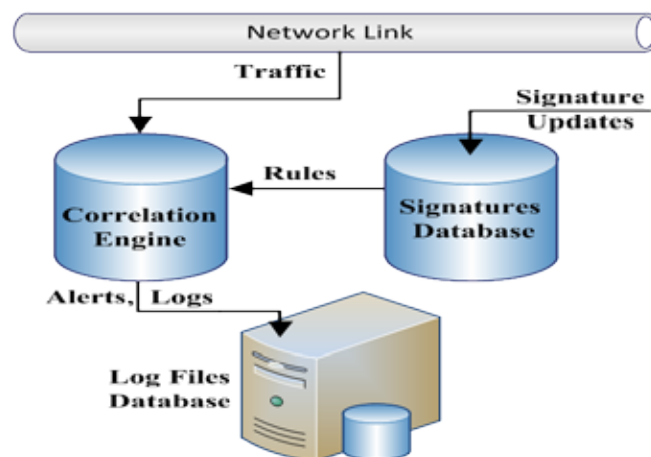


Figure 2.9: a simplified figure representing a signature-based detection model