

AN EMPIRICAL STUDY OF AN ANONYMITY METRIC FOR DATA NETWORKS

by

Abinash Vasudevan
A Thesis
Submitted to the
Graduate Faculty
of
George Mason University
in Partial Fulfillment of
The Requirements for the Degree
of
Master of Science
Electrical Engineering

Committee:

Brian L. Mark

Dr. Brian L. Mark, Thesis Director

Yariv Ephraim

Dr. Yariv Ephraim, Committee Member

KGaj

Dr. Kris Gaj, Committee Member

Andre Manitus

Dr. Andre Manitus, Department Chair

Kenneth S. Ball

Dr. Kenneth S. Ball, Dean, Volgenau School
of Engineering

Date: 08/23/2012

Fall Semester 2012
George Mason University
Fairfax, VA

An Empirical Study of an Anonymity Metric for Data Networks

A thesis submitted in partial fulfillment of the requirements for the degree of Master of Science at George Mason University

By

Abinash Vasudevan
Bachelor of Engineering
Saveetha Engineering College, 2009

Director: Brian L. Mark, Professor
Department of Electrical and Computer Engineering

Fall Semester 2012
George Mason University
Fairfax, VA

Copyright 2012© Abinash Vasudevan
All Rights Reserved

Dedication

This is dedicated to my parents.

Acknowledgements

I would like to take this opportunity to express my sincere gratitude to Dr. Brian Mark. Without his patient guidance and support, none of this is possible. I would also like to thank my parents for their continuous support and trust.

Table of Contents

	Page
List of Figures.....	vii
List of Abbreviations	viii
Abstract.....	xi
1. Introduction.....	1
1.1 Background.....	1
1.2 Problem Statement.....	3
1.3 Organization.....	4
2. Review of Anonymity Metric.....	5
2.1 Information Theoretic Anonymity Metric.....	6
2.1.1 Knowledge Vulnerabilities.....	7
2.1.2 Entropy.....	8
2.1.3 Route Length.....	10
2.2 Normalized Information Theoretic Metric.....	11
3. Background.....	15
3.1 Timing Analysis.....	15
3.1.1 Timing Analysis in a Mix-Based System.....	16
3.1.1.1 Attack Model.....	17
3.1.2 Traffic Pattern Analysis in TOR based Anonymous Network.....	19
3.1.2.1 Attack Model.....	20
3.2 Traffic Models.....	22
3.2.1 Poisson Process.....	22
3.2.2 Markov Modulated Poisson Process.....	24
3.2.3 Markov Arrival Process.....	26
4. Wavelet-Based Anonymity Metric.....	29
4.1 Wavelet Transformation.....	29
4.1.1 Haar Wavelet.....	31
4.2 Wavelet Based Anonymity Metric.....	31
4.2.1 Network Assumptions.....	32
4.2.2 Timing Distortion.....	34
4.2.3 Multi-Resolution Analysis.....	36
4.2.4 Properties of Energy Based Metric.....	38

5. Experimental results.....	40
5.1 Energy Plots for Poisson Process.....	41
5.2 Energy Plots for Different Inter-Arrival Models.....	47
5.2.1 Theoretical Calculation of Arrival Rate for MMPP AND MAP.....	48
5.3 Energy Plots for MMPP and MAP traffic Arrival.....	50
5.4 Conclusion.....	52
6. Conclusion.....	54
Bibliography.....	56

List of Figures

Figure	Page
2.1 Example network for Knowledge Vulnerabilities.....	8
2.2 Using maximum route length to reduce the Anonymity size.....	11
3.1 Mix-Based System.....	17
3.2 Timing analysis using two paths.....	18
3.3 Generalized Poisson process.....	24
3.4 Poisson process with same arrival rate.....	24
3.5 MMPP state transition diagram.....	26
3.6 MAP state transition diagram.....	28
4.1 Haar Wavelet.....	31
4.2 Low Latency Anonymous Network.....	33
4.3 Packet Timing Distortion.....	35
5.1 Energy plot for 3 different Poisson Rates.....	42
5.2 Energy plot for different delays for inter-arrival time of 0.9 seconds.....	44
5.3 Energy plot for different delays for inter-arrival time of 1.5 seconds.....	45
5.4 Energy plot for different delays for inter-arrival time of 2.1 seconds.....	46
5.5 Energy plot for different distributions of inter-arrival time.....	47
5.6 Energy plot for different delays for MMPP inter-arrival time.....	50
5.7 Energy plot for different delays for MAP inter-arrival time.....	51

List of Abbreviations

1. MRA - Multi-Resolution Analysis
2. MMPP - Markov Modulated Poisson process
3. MAP - Markovian Arrival Process

Abstract

AN EMPIRICAL STUDY OF AN ANONYMITY METRIC FOR DATA NETWORKS

Abinash Vasudevan, M.S.

George Mason University, 2012

Thesis Director: Dr. Brian L. Mark

Privacy and data protection are two very important needs in the modern day Internet. One of the attacks to privacy is eavesdropping, i.e., an outsider or an attacker listens to a private conversation and identifies the people involved in the conversation. In order to protect the content which is to be transmitted, encryption methods are used. Even if the data is encrypted, it is possible for the attacker to identify the end user, i.e., the person sending the data and the person to whom the data is sent.

An anonymous network is a type of network that prevents traffic analysis and protects the identity of the end users. Some of the popular anonymous networks, namely Anonymizer.com, TOR, etc., are used for identity concealment. These anonymous networks provide real-time, low latency anonymous communications. Because of the low latency implementation, timing constraints are imposed on the low-latency networks. Due to this timing constraint, an attacker will be able to get details from the packet timing

information and use it to identify the end users. This makes the timing attack possible in these low-latency anonymous networks.

In this thesis, an anonymity metric that can measure the practical effectiveness of low-latency anonymous network is studied. The anonymity metric calculates the timing distortion between the incoming packets and the outgoing packets and uses wavelet-based Multi-Resolution Analysis (MRA) to determine the anonymity of the network. For the purpose of analysis, packet traffic is simulated based on several stochastic traffic models. The simulated traffic contains information about the timing of packets entering and leaving the anonymous network. The most basic traffic model used for this purpose is the Poisson process; we also use the Markov Modulated Poisson Process and the Markovian Arrival Process for further analysis. The end-to-end network delay is characterized with various probability distributions for this study. The results of the analysis show that the measured energy used to compute the anonymity metric is higher for more complex traffic arrival models, which implies that the anonymity level is correspondingly higher. The more complex traffic patterns introduce more randomness in the timing information, resulting in higher measured energy values. Hence, for example, the Poisson process arrival model has the least energy among the traffic models used in this study.

Chapter 1: Introduction

1.1 Background

An anonymous network provides a mechanism for concealing the identity of the sender and the receiver. In the modern day Internet, privacy is a very important issue. Traffic analysis is a method that can be used to reveal the identity of the end users, which the user may not want to disclose. For example, e-mail users may not want others to know with whom they are communicating as well as their own identities. Users wanting to use the Internet for financial transaction may want to remain anonymous. There are a few anonymizing techniques which can protect the identity of both the sender and receiver. Some of these include onion routing [1] and mix networks [2].

Onion routing [1] is a routing technique that can protect the identity of the end users. It uses proxies and data encryption to achieve anonymity. As the name implies, the data is encrypted multiple times, forming layers above layers, like an onion. If user A wants to send a message to user B, the onion routing protocol determines a path in a network made of so-called onion routers and constructs a layered data structure called an Onion.

A mix network [2] is another type of anonymizing technique to protect the identity of people involved in the communication. A mix is a proxy that tries to conceal the relationship, mainly timing information between the incoming packets and the outgoing

packets. A mix uses a wide range of methods to achieve this level of anonymity for the users. It will reorder the packets, delay the packet, and also use encryption. Along with these the mix network will add chaff traffic, split the traffic, merge the traffic, and perform packetization and repacketization. Adding chaff traffic is the process of inserting meaningless padding packets to the original traffic. All of these operations are done to basically eliminate any kind of correlation from the timing information so that it will be hard for the attacker to correlate the timing information for all packets.

Traffic analysis is done to identify the sender and receiver. There are some methods that can be used to do this and one of them is timing analysis [3]. For a high latency network or application which is not time-constrained, timing analysis will be very hard to carry out. Because it is not time bound, it can take a long time for the message or data to be transferred. During this time it can go through many transformations and packet manipulation techniques before reaching the destination, thus making it hard to do traffic analysis. But when it comes to real-time and low-latency networks, timing analysis is very much possible because of the timing constraints. In timing analysis, the attacker studies the time stamps of the packets sent or received through the system and calculates the correlation between them. There are a few timing analysis methods described in the literature references [3, 6, and 7] which expand about the problem with low latency anonymous networks.

An anonymity metric is used to characterize quantitatively the level of anonymity provided by an anonymous network. All the anonymous networks require some form of validation to assess how successfully the identity of the end users is concealed.

Calculating an anonymity metric will give a basic understanding of the performance of a given anonymous network and help in future development of the network.

1.2 Problem Statement

Anonymous networks for low-latency connection and real-time application are susceptible to timing attacks. Most of the work relating to anonymous networks is concerned with identifying the problems associated with a specific method, e.g. traffic analysis method, watermarking techniques, etc. But calculating the effectiveness of an anonymous network has not been done very extensively. This is very important because there should be some form of quantitative approach to characterize the anonymity level of a particular anonymous network. Also, it can be used to compare different kind of anonymous networks. Measuring anonymity for each network is a more reasonable way to compare them. If a new method is proposed for anonymous communications, an anonymity metric can be used to test it or any updates that are done for an already existing network. This enables the network provider to check whether any actual improvement was achieved in a network upgrade.

In this thesis, a novel method of quantifying the practical effectiveness of anonymous networks proposed in [12] is studied empirically. This method uses the distortion in the timing information of the packets sent through the system. The timing information is analyzed using wavelet-based multi-resolution analysis. This analysis provides an anonymity metric for a particular network. The anonymity metric that is measured can be applied to any kind of anonymous network, since it is based only on timing information.

1.3 Organization

The organization of this thesis from here on is as follows. Chapter 2 provides a survey of several anonymity measurement techniques that have been proposed in the literature.

Chapter 3 contains information about the traffic analysis methods and a description of stochastic traffic models used for the simulation study. In chapter 4, we discuss the wavelet-based anonymity metric proposed in [12]. In chapter 5, we present and discuss the simulation results. Chapter 6 concludes the thesis.

Chapter 2: Review of Anonymity Metric

Anonymity, in general, refers to the concealment of one's identity. When it comes to remaining anonymous in the Internet, for example, a person or a computer, whatever role the user may be (receiver or sender) in association with a data packet should not be identifiable. If that user cannot be identified then that user is said to be anonymous. In real world scenarios, it is very hard for a user to remain completely "anonymous." But relative anonymity can be achieved with the help of specially designed network architectures. Relative anonymity has a slightly different meaning than "absolute" anonymity. Consider, for example, a data packet is sent by a user and there are 3 possible senders A, B and C who could have sent the data. The attacker in this case may know the existence of all the possible senders but may not be able to find out the real sender. The attacker will know that there are three users that could be the sender. So for this specific case if the attacker sees all the users as equally likely to be the sender, then relative anonymity is achieved i.e., the attacker cannot find the actual sender from the possible senders.

Listing the properties of anonymity is more difficult. It depends on the needs of that specific user as well as the parameters used to find anonymity. Also, on the other hand, the users would like to know the amount of anonymity being offered by an anonymous network to them. Hence, assessing anonymity that is offered by any anonymous network

is very important. For this purpose, several techniques have been proposed in the literature that can be used to evaluate the anonymity that is provided by a network. Each of the techniques has its own advantages and disadvantages.

2.1 Information Theoretic Anonymity Metric

Concept from information theory can be used to determine the anonymity level based on the amount of information lost. In [4], a method that can be used to calculate the amount of anonymity offered in a network is discussed. This paper analyzes the anonymity of a message going through a mix-based anonymity system. In the mix-based system, nodes called mixes that are used to create anonymous connections. In such a system, the sender, instead of passing the message directly to the recipient forwards it through a number of mixes.

Various definitions for anonymity have been proposed in the literature.

The definition for anonymity given in [4] is as follows: “Anonymity is the state of being not identifiable within a set of subjects, the anonymity set.” The concept of anonymity and anonymity are refined further as follows: “Anonymity is stronger, the larger the receptive anonymity set is and the more evenly distributed the sending or receiving, respectively, of the subjects within that set is.”

This even distribution of sending or receiving of members forms a new requirement for characterizing the degree of anonymity provided by the network. Therefore using only the size of the anonymity set is not a sufficiently good representation of anonymity. It will not form the complete requirement for calculating anonymity in a network and will

not show the complete picture of the actual anonymity level achieved by the network. There are more things other than the size of the anonymity set to characterize the anonymity level of a particular network. For example, one user from a particular anonymity set may send more data compared to other users. More data here can mean a greater number of packets as well as more time spent on the data channel. In this case, that particular user is more likely to send data compared to other. So with time, the attacker will come to know of this pattern and use this information to identify different users who are connected to a particular communication stream.

2.1.1 Knowledge Vulnerabilities

There are a few vulnerabilities that are associated with the usage of the anonymity set metric in the presence of an attacker with additional information. The additional information that the attacker has, can be from a compromised mix or with respect to a traffic pattern with a particular anonymity set. A mix is an intermediate node in the communication link that receives the packet in the anonymous network and sends it to another mix. A compromised mix is a node in the network which is controlled by an attacker. The ultimate aim of the anonymous network is to reach the destination from the source through a number of mixes without revealing the end users' identities.

Consider the arrangement of mixes as in Figure 1. The small squares in the Figure represent senders, labeled as A, B, C, D, E and F. The bigger boxes are mixes with a threshold of 3 (each mix can have a maximum of 3 input lines and 3 output lines). Some of the receivers are named with their sender anonymity set. Suppose that there are two

users sending data packets and two users receiving data packets. In this setup, if the attacker somehow establishes the fact that for example, A is communicating with V, he can derive the fact that F is communicating with Z. So in order to establish a link between F and Z, all the attacker has to know is that one of A, B, C, D, E is communicating with V. This kind of information is not included if the anonymity set size is used directly as a measure of anonymity. Also, not all receivers or senders suffer from a problem like this. So the pattern of the traffic or the connectivity of the network may lead to various kinds of information leakage.

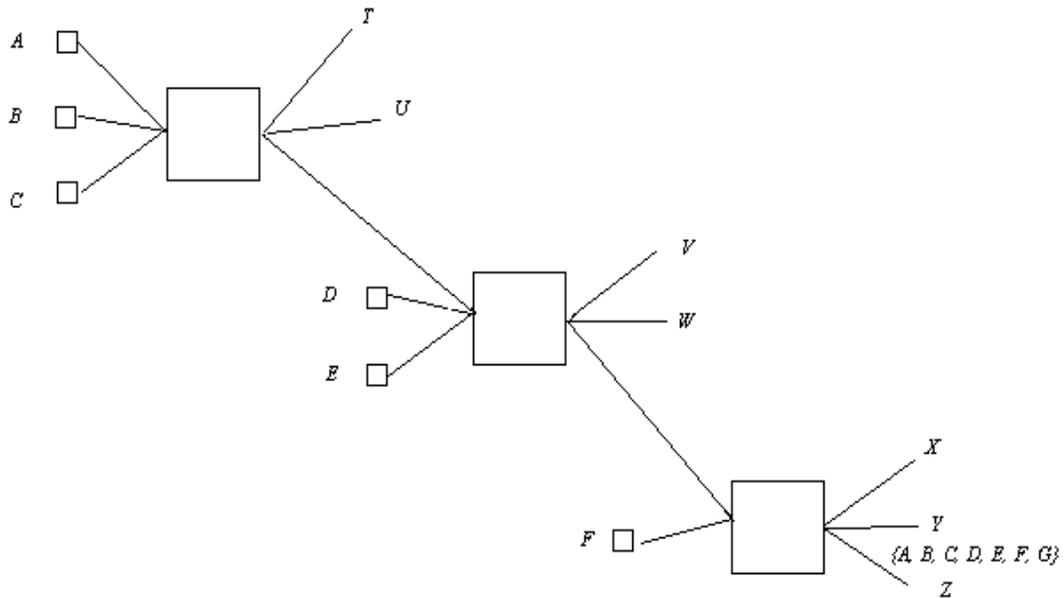


Fig.2.1 Example network for knowledge vulnerabilities

2.1.2 Entropy

As seen from above discussion it is clear that, the size of the anonymity set alone cannot be used to calculate the anonymity level in a network. The concept of entropy is used for obtaining the probability distributions of a user for being a sender or a receiver with respect to the role played by it in message M . An anonymous communication model is defined as follows [4]:-

“Given a model of attacker and a finite set of all users ψ , let $r \in R$ be a role for the user ($R = \{\text{sender, recipient}\}$) with respect to a message M . Let U be the attacker’s a-posteriori probability distribution of user $u \in \psi$ having a role r with respect to M .”

This model does not include the size of the anonymity set. Instead, it uses an anonymity probability distribution U . Given a message m , we have a probability distribution of its possible senders and receivers as viewed by the attacker. The attacker may assign zero probability to some users. It means that the specific user did not have any role to play with respect to message M . On the contrary, the attacker can assign a maximum probability of 1 to any one user which means, that specific user is the only user associated with the message M for role r . For example, if there is an attacker with the network setup as in the mix system of Figure 1 and the message considered is seen by the attacker as having arrived at T , then $U(\text{receiver}, T) = 1$ and $\forall S \neq T U(\text{receiver}, S) = 0$. If all the users are assigned an equal probability, then the size of the anonymity set can be used to characterize the anonymity level. This equiprobable case occurs when the

network achieves the highest anonymity level. So the effective size of the set with different non-zero probabilities for each user is

$$S_{eff} = - \sum_{u \in \psi} p_u \log_2 p_u$$

where $p_u = U(u, r)$.

This effective size S_{eff} is equal to the entropy of the distribution. This can be interpreted as the additional information that is needed by the attacker to definitively identify the user u with role r for a particular message M .

2.1.3 Route Length

With the probabilities included in the model, the effective size of the system will be different than the original anonymity size. The attacker will be trying to reduce the effective size further in order to identify the user. Route length plays an important role to this reduction of the effective size. Some arrangements of the system make it more vulnerable to route length and narrowing the anonymity set is achieved easily. In a mix based system, the number of mixes that the message passes through is fixed when the message is created. Hence, say if a compromised mix comes to know about information like the number of mixes the message went through till now, it can calculate how many mixes it still needs to reach the destination. This information will strengthen the attack on the anonymous network. This problem is illustrated with an example. Consider the situation in Figure 2, where each arrow represents a message observed by the attacker.

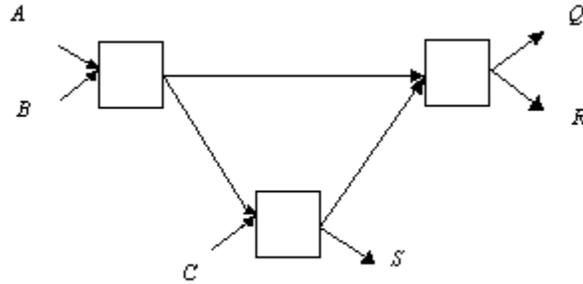


Fig.2.2 Using maximum route length to reduce the Anonymity size.

Let the maximum route length be assumed to be 2, i.e. any message can pass through 2 mixes and no more than that. So the arrow from the bottom mix to the rightmost mix can have only C as the sender. If the message is from A or B it would take 3 mixes to reach the destination. Hence, C is the only possible sender. So the new anonymity set would be {A, B} with the inclusion of maximum route length. Otherwise it would have been {A, B, C}.

2.2 Normalized Information Theoretic Metric

Another method for measuring the anonymity of a network using a similar method is proposed in [5]. Information theory is used in calculating the anonymity that the users get in a network, similar to [4]. This paper focuses on connection anonymity and does not focus on data anonymity. There are a few differences between data anonymity and connection anonymity. Data anonymity is about identifying the information that is exchanged in a particular application. Connection anonymity is about hiding the identities

of the source and destination during the data transfer. The paper [5] proposes a model, based on Shannon's definition of entropy that allows quantifying the degree of anonymity in a network. The degree measured will be dependent on the attacker's ability. This paper considers the system to provide anonymity through mixes. The system model consists of entities like Senders, Recipients and Mixes. Senders are users who send messages to recipients and can be grouped into a set called the anonymity set. The messages can be emails, requests to web servers or any data stream. Recipients are entities that receive messages from the senders. Recipients can be active (if they reply back to senders) or passive (if they do not reply back to senders). There are large varieties of recipients like web servers, email accounts and so on. The last entity that forms the system is mixes. Mixes are the nodes that are used to create anonymous connections. They take messages and send them out so that correlation between them is decreased.

The attack model depends on the probabilities that a particular user is the sender. Also, the degree of anonymity is measured with respect to a particular message and the network setup. So if there is any change in network setup, the attacker has to repeat all the steps from the beginning which makes the existing attack analysis not valid. The reason for this is that, when the network setup is changed (like one of the users leaves the network or a new user is added) there might be a change in the pattern and hence, the probabilities assigned to each user should be changed.

The definition in [5] for anonymity is the state of being not identifiable within a set of subjects called the anonymity set. The anonymity set in this case is the set of honest users who might send a message. So in a system of N users, maximum anonymity is achieved,

if the attacker sees all the users as equally likely to be the originator of the message. The method that is used in [5] depends on the probability distribution of the users that they are the source of the message. By this way, the measurement of anonymity is made less dependent on the size of the anonymity set. The model compares the information obtained by the attacker after observing the system against the optimal situation, in which all honest users seem to be equally probable as being the originator of the message. In a system with N users, then each user would be assigned a probability of $1/N$. After observing the system for a while, the attacker assigns probabilities to each user as being the sender of the message. This is based on the information leaked by the system during the course of observation by means of timing attacks, timing analysis and more.

Let X be a discrete random variable with probability mass function $p_i = Pr(X = i)$, where i represents each possible value that X may take. In this case each i corresponds to an element in the anonymity set. $H(X)$ represents the entropy of the system after the attack has taken place. So to each sender belonging to the anonymity set of size N , the attacker assigns a probability p_i and $H(X)$ is calculated as follows

$$H(X) = - \sum_{i=1}^N p_i \log_2 p_i$$

Let H_M be the maximum entropy of the system with N users in the anonymity set.

$$H_M = \log_2 N ,$$

where N is the total number of honest sender.

The information that the user got with the attack can be expressed as $H_M - H(X)$. The degree of anonymity d that is provided is defined as

$$d = 1 - \frac{H_M - H(X)}{H_M} = \frac{H(X)}{H_M}$$

It is divided by H_M to normalize the value. $H(X)$ can be considered as the additional information that the attacker needs to identify the sender. For any case, the value of d will be in between 0 and 1, i.e., $0 \leq d \leq 1$. In a particular system, a user or a small group of users are shown as originators with a high probability with respect to others, then the system is not providing a high degree of anonymity. As said before, any system with equiprobable distribution will provide a degree of anonymity of one. For example, in a system of 2 senders, the degree of anonymity $d = 1$ if both of them are assigned probability of $1/2$. The information learned by the attacker during the attack is $H_M - H(X)$. One minus this term gives the anonymity level achieved by the network. Hence, the degree of anonymity provided by the system quantifies the amount of information leaked by the system during observation.

Chapter 3 Background

3.1 Timing Analysis

Timing analysis is a method that is used by an attacker to disclose the identity of the users, who are communicating. Timing analysis studies the timing stamps of messages moving through the system to find correlations. Today, the Internet is widely used for real-time communications like web browsing, instant messaging and the demand for this type of communication is ever increasing. With the advancement in technology, the latency between the communicating users is reduced more and more. This makes it harder to achieve anonymity since it takes less time to transfer the packets to the end user.

The effect of the timing attack in a low latency network is much higher than in a high latency network. The applications that use a low latency network are generally real-time applications which require less time for communication between the users. Hence, the application that is responsible for communication dispatches a packet as and when it is formed (instantaneous dispatch to reduce the delay). Due to the time constraints, the packets cannot be manipulated in too many ways. Hence, the information that is leaked by a low latency network is higher than that of a high latency network. In a high latency network end-to-end delay is not a problem. Since delay is not a problem, a packet can be manipulated in many ways in a high latency network. Two of the methods of

manipulation are 1) the packets can be delayed at some point in the communication link and 2) the packets can be reordered. These methods make it hard for the attacker to correlate the timing information. So there is an inherent problem that is associated with low latency networks when it comes to anonymity. In this section we provide some background about timing analysis.

3.1.1 Timing analysis in a Mix-Based System

A mix system, also called simply a “mix” is a communication proxy that helps to hide the correspondence between its incoming messages and outgoing messages [6]. The communication in a mix network happens through a chain of mixes which provides unlinkability of senders and receivers despite the presence of an observer and compromised mixes. A typical mix will be able to delay packets, reorder them or even emit additional dummy packets along with the original packets. This makes it very effective for a high latency network to create anonymity but not for a low latency network. The attacker studies the timings of the messages through the system to find correlations between them. The kind of analysis that is made by this method is a two-attacker mixes method [6] (i.e., compromised mixes or owned mixes by the adversary) to determine the communicating users. The two-attacker mixes method is one which uses two or more compromised mixes to identify the path used by the sender and the receiver. Attacker mixes are the mixes that are controlled by the attacker.

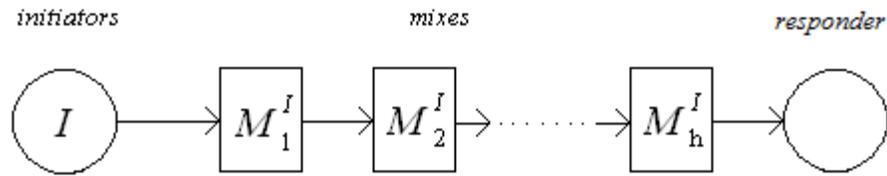


Fig 3.1 Mix-based System

Figure 3 illustrates a communication path in a mix system between an initiator and a responder. The initiator I is one who sends the packet (source) and a responder is one which acts as the other end of the communication. The initiator I uses a path P^I , of the mixes of the system. Path P^I consists of h mixes that start from M_1^I and ends with M_h^I . Mix M_1^I is the first mix in the path that is connected to the initiator. For convenience it is assumed that mix h is the last mix in the chain of mixes, i.e., Mix M_h^I is the last mix in the path and connects to the responder. Mix M_1^I receives the packet from the initiator and mix M_h^I sends the packet to the appropriate sender. We assume that each link between the two mixes typically carries packets from multiple initiators, and that each packet received, a mix can identify the path P^I to which the packet corresponds. Also, M_1^I will know it is the first mix in the path and M_h^I will know it is the last mix in the path.

3.1.1.1 Attack Model

The ultimate aim of timing analysis is to find a correlation between the timings of packets seen by M_1^I and those seen by an end point M_h^I . Packets that initiator I , sends along the path P^I runs along general purpose link, between the initiator and the first mix as well as

between mixes. These general purpose links between mixes will have delays or packet drops associated with it. Each packet will be affected by these delays or drops and forms the basis on which the attacker finds correlation. In this method the attacker controls M_1^I and M_h^J on two paths P^I and P^J and the goal of the attacker is to determine whether $I = J$. The stronger the correlation between M_1^I and M_h^J , the more likely $I = J$ and M_h^I is actually M_h^J . There is a possibility that M_1^I and M_1^J can see the same timings of packets and then it will be very difficult to match the packet stream with M_h^J .

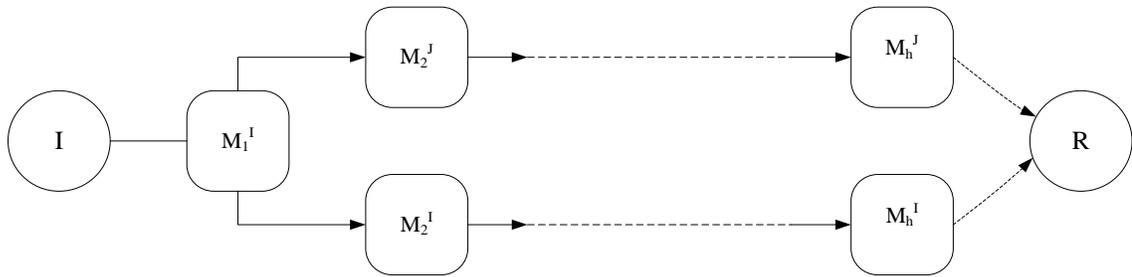


Fig 3.2 Timing analysis using two paths

To study the timing correlation, a random variable is defined. The timing difference δ_i is the difference in time between the arrival of packet i and the arrival of its successor packet. If both the mixes are on the same path P^I , i.e., if M_1^I and M_h^I is in path P^I , it will have a high correlation. The problem with this variable is that it is very sensitive to dropped packets. A dropped packet between the first mix and the last mix in the path will cause the timings to be off by one. This makes the correlation that is calculated between the mixes not a match. To avoid this problem to some extent, a new

random variable is introduced that is obtained from the data. It uses non-overlapping and adjacent windows of time with fixed duration W . In each window k , mix M maintains a count X_k^l of the number of packet arrivals on the path P^l and hence, reduces the effects of packet drop. With an active attacker, the effects of timing analysis are even more.

3.1.2 Traffic Pattern Analysis in TOR Based Anonymous Network

Timing attacks give lot of information that the attacker can use to identify users. Similar to the previous method, [7] uses the timing information for find the identity of users. It uses a TOR network to test the method of the timing attack. TOR [8] is type of anonymous network that is popularly used for low latency anonymous communication. The TOR network is composed of set of nodes that act as relays for a number of communication streams with different users. A stream is a communication line between two different users using TOR network. The most important function of a TOR node is to hide any correlation between incoming packets and outgoing packets. This will make it difficult for the attacker to find the actual users who are communicating using the TOR network. TOR uses onion routing for end to end communication between users. Onion routing performs multiple layers of encryption on the original data packet and forwards it. Each time this packet reaches an onion router, a layer of encryption is stripped and then forwarded to the next onion router. By this way, the attacker cannot correlate packets based on the content.

The TOR network is formed using a set of nodes. This acts as a circuit for a particular stream for transfer of data packets. During this circuit formation, secret keys

are negotiated with each TOR node and this establishes a secure channel. All the communication between the users is tunneled through this channel along with multiple layers of encryption for the data packets. Hence, whenever a packet reaches a TOR node, it strips the top encrypted layer and then forwards it to the next onion router. This is done till it reaches the final TOR node beyond which it is the original data for the intended responder without any encryption. Unlike mix type anonymizer techniques, TOR does not mix or modify data packets explicitly. Each data stream will have a separate buffer and when packets arrive for a specific stream, it stores it in that specific buffer allotted for that particular stream. TOR now uses a round-robin method to dispatch the packets stored in buffer. If a connection or stream does not have any packet stored in its buffer, then it skips that stream and goes to the next non-empty buffer in the sequence. TOR is predominantly used for communication which requires less end to end delay. Due to this reason, TOR nodes cannot perform actions that modify the packet flow like delaying, reordering or dropping packets. The attacker uses this property to exploit the timing information of the packets and tries to identify the initiator and the responder.

3.1.2.1 Attack Model

The ultimate objective of attacking an anonymous communication network is to link the initiator with the respective responder. As with other few anonymous communication models, TOR assumes a weak threat model as a reference. It can protect the identity from a non-global adversary who can control or see only a small portion of the network. This particular attack model used in [7] uses the traffic pattern in a particular node to identify

the other TOR nodes that stream is connected and ultimately the users (sender and receiver).

Unlike conventional traffic analysis that uses time stamps for identifying the end users, [7] uses a slightly different method. The ability to use a TOR network and route over the anonymous communication network can be used to estimate the traffic load on a specific node. As discussed before, each TOR node uses a round robin method to dispatch packets in each stream. So the latency for each packet is determined by the load in each node. This means if there is a higher load in a specific TOR node due to one particular stream, it will affect the latency characteristics of all the packets through other communicating streams. The attacker may control one or more TOR nodes in the network (non-global attacker). These corrupt nodes form a connection with other TOR nodes that are to be analyzed. Each node based on the load has a different traffic pattern. The corrupt TOR node analyzes the pattern and compares it with a known traffic pattern and determines whether that node is used for a particular communication. If the attacker controls a web server on top of TOR nodes, the results will reflect more accurate measurements. The attacker-controlled web server sends data packet to the communication initiator in a specific pattern. Now the corrupt TOR node checks for this pattern in each and every node it is connected to. By this way it can back track to the initiator and ultimately reveal the identity.

3.2 Traffic Models

Timing information forms the basis for timing analysis and the attack, used to identify the end users. This means that the time stamps of packets are used, i.e., the time at which a particular packet enters the anonymous network and the time the packet leaves the anonymous network. Instead of getting data (timing information) from a real-time anonymous network, they can be simulated. To generate this timing information (inter-arrival time between packets, rate of packet arrivals etc.) for simulation, stochastic models can be used. Depending upon the distribution, the characteristics of the timing information varies. The random processes that are of particular interest are the Poisson process, Markov Modulated Poisson Process (MMPP) and Markovian arrival process (MAP). The following sections provide a brief introduction to these random process models.

3.2.1 Poisson Process

In probability theory and statistics, the Poisson distribution expresses the probability of a given number of events occurring in a fixed interval of time. An arrival or an event can be a simple occurrence of certain things in time like, for example, an arrival of phone calls to a switch board or an arrival of e-mail messages.

Let X be a Poisson random variable, P denote the probability measure and let Ω be the sample space. Let λ be the mean arrival rate of the Poisson random variable. The probability of k arrivals in an interval of length T is defined as

$$P(X = k) = \frac{\lambda T^k e^{-\lambda T}}{k!}$$

This gives the probability of k arrivals in the specified time period. The inter-arrival times between arrivals are independent and distributed exponentially with mean parameter $1/\lambda$.

The Poisson process is a continuous time process with the inter-arrival time following exponential distribution. Let Ω be a sample space and P a probability measure on it. An arrival process $N = \{ N_t ; t \geq 0 \}$ defined on Ω such that for any $\omega \in \Omega$, the mapping $t \rightarrow N_t(\omega)$ is non-decreasing, increases by jumps only, is right continuous, and has $N_0(\omega) = 0$. An arrival process $N = \{ N_t ; t \geq 0 \}$ is called a Poisson process if it satisfies the following axioms [9]:

- (a) For almost all ω , each jump $t \rightarrow N_t$ is of unit magnitude;
- (b) For any $t, s \geq 0$, $N_{t+s} - N_t$ is independent of $\{ N_u ; u \leq t \}$;
- (c) For any $t, s \geq 0$, the distribution of $N_{t+s} - N_t$ is independent of t .

In the definition, axiom (b) shows independence of number of arrivals in $(t, t + s]$ from the past history until t . A generalized Poisson process may be viewed as a Markov process with generator matrix given by

$$Q = \begin{bmatrix} -\lambda_0 & \lambda_0 & 0 & \dots & \dots & 0 \\ 0 & -\lambda_1 & \lambda_1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \ddots & \ddots & \ddots & \ddots \end{bmatrix}$$

This is a generator matrix where each λ_i defines the arrival rate in each state.

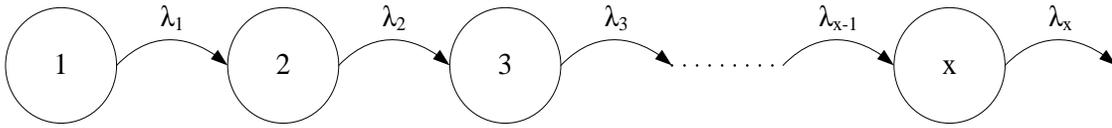


Fig.3.3 Generalized Poisson process

In the homogenous case all λ_i 's are equal to λ .

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 & \dots & \dots & 0 \\ 0 & -\lambda & \lambda & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \ddots & \ddots & \ddots & \ddots \end{bmatrix}$$

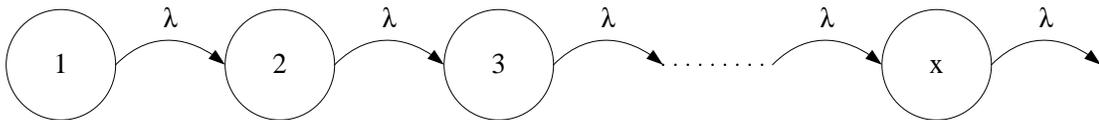


Fig. 3.4 Poisson process with same arrival rate

3.2.2 Markov Modulated Poisson Process

Markov Modulated Poisson Process (MMPP) is a doubly stochastic Poisson process with an irreducible underlying Markov process. The MMPP [10] is a generalization of Poisson process. An MMPP can be constructed by varying the arrival rate of a Poisson process according to m-state irreducible continuous time Markov chain which is independent of the arrival process. When the Markov chain or the process is in state i , then it is characterized with a Poisson arrival rate of λ_i . The generator matrix Q for an MMPP is

$$Q = \begin{bmatrix} -\sigma_1 & \sigma_{12} & \sigma_{13} & \cdots & \sigma_{1\ m-1} & \sigma_{1m} \\ \sigma_{21} & -\sigma_2 & \sigma_{23} & \cdots & \sigma_{2\ m-1} & \sigma_{2m} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \sigma_{m3} & \cdots & \sigma_{m\ m-1} & \sigma_m \end{bmatrix}$$

Where,

$$\sigma_i = \sum_{\substack{j=1 \\ j \neq i}}^m \sigma_{ij}$$

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_m),$$

$$\lambda = (\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_m)^T$$

For homogenous Markov Modulated Poisson process Q and Λ are independent of time.

The steady-state vector of Markov chain π is defined as follows,

$$\pi Q = 0,$$

$$\pi e = 1,$$

where $e = (1, 1, 1 \dots \dots 1)^T$ is the column vector of length m .

The state transition diagram will give a much clear picture of what Markov Modulated Poisson process does. In Figure.5, the state-transition diagram for an MMPP is given. As seen from the diagram there are m states in the Markov model. The model is characterized by two random process S and N . S here shows the current state it is in and N is the counting process (counts the number of arrivals). From each state i it can move to any other state with same count at the rate of σ_{ix} , x being the new state such that, $1 \leq x \leq m$ and $x \neq i$. The count N can increase by 1 with each arrival, with arrival rate of λ_i . The process starts with $(1, 0)$ state. From here it can go to any state x $(x, 0)$ with count 0 or it can go to $(1, 1)$ state with an arrival.

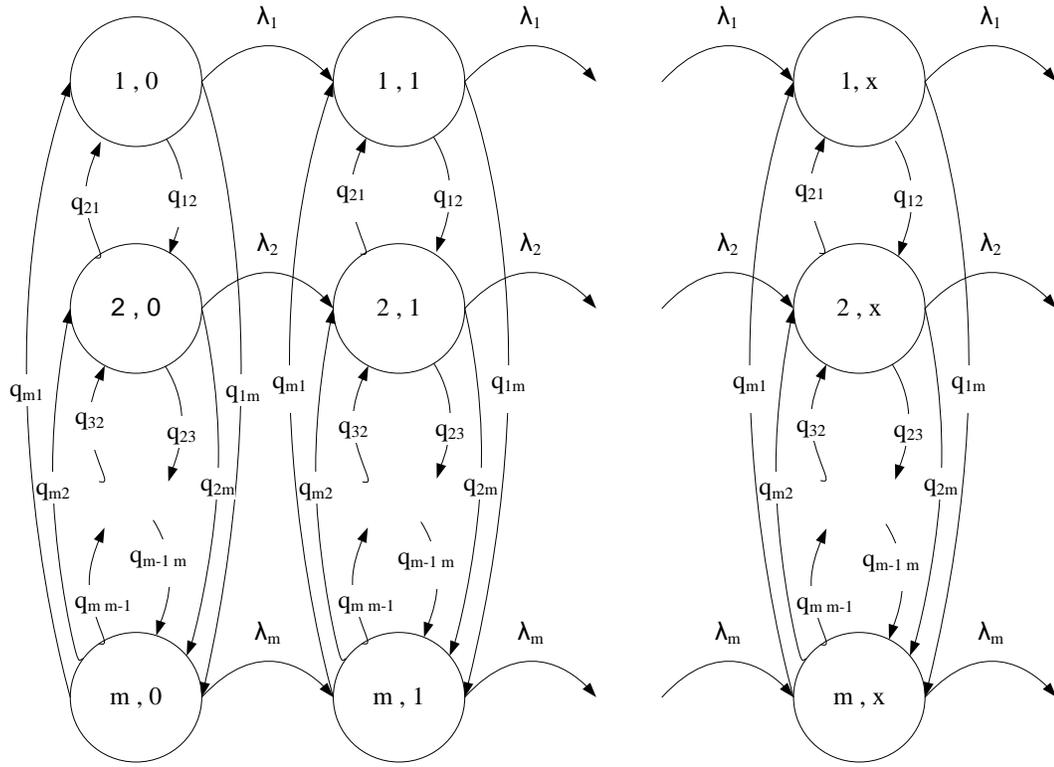


Fig. 3.5 MMPP state transition diagram

3.2.3 Markovian Arrival Process

The Markovian Arrival process (MAP) [11] is a generalization of the MMPP. Like the MMPP it is represented by two processes, S and N . S denotes the current state and N denotes the counting process. Like the MMPP, the count N can be incremented by only one not more than that. So it can go to a different state of same count or to any state of the next incremented count. This is the most basic difference between the MAP and the MMPP. The MAP increments the count by one and stays in that state for a non-exponential time period and maintains the underlying Markov structure.

$$Q = \begin{bmatrix} D_0 & D_1 & 0 & \dots & 0 & \dots \\ 0 & D_0 & D_1 & \dots & 0 & \dots \\ 0 & 0 & D_0 & D_1 & \dots & 0 \\ \dots & \dots & \dots & \ddots & \ddots & \ddots \end{bmatrix}$$

The Q matrix is called the generator matrix. D_0 and D_1 are $m \times m$ matrices. D_0 denotes a transition matrix for no arrival (S changes and N stays same). D_1 denotes the transition matrix for an arrival (both S and N can change; N can change with increments of 1).

After a non-exponential sojourn time, the process jumps from state (j, i) to state $(l, i+k)$ where $k = 1$ or $k = 0$, i.e., either there is an arrival or not and l can be any underlying state for the event arrival.

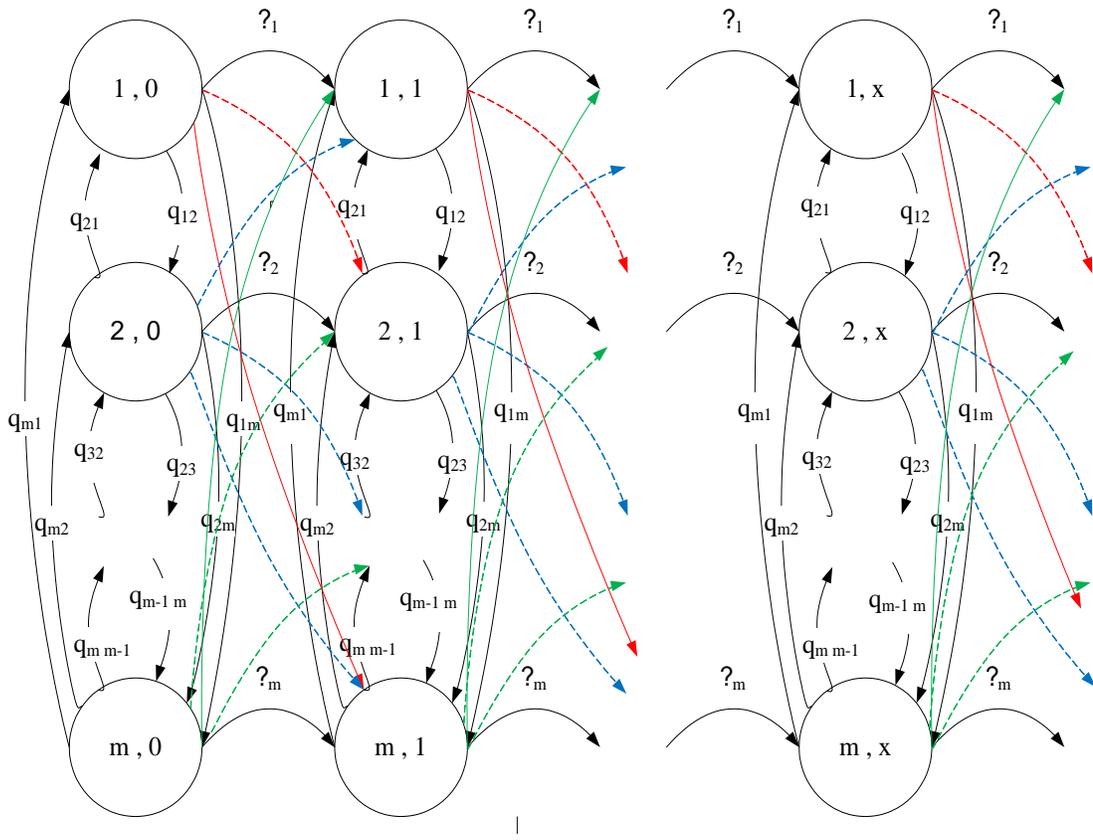


Fig 3.6 MAP State transition diagram

Chapter 4: Wavelet-Based Anonymity Metric

As discussed in chapter 3, the timing attack causes serious problems to the anonymity in an anonymous network. In a low latency network, the impact of a timing attack is many times higher than in a high latency network. The number of Internet users is growing steadily and most of them use real-time web communications like browsing, chatting etc. The end-to-end delay must be low for real-time communication, which is consequently more prone to timing attacks. Due to these factors, protecting the identity of the users is very difficult in the presence of timing attacks. One other problem with the anonymity metric is measuring it. Each network might be unique in its arrangement and how different users are connected between them. So the method that is used to calculate the anonymity level should be designed in a way such that it should be independent of the network or connections. This chapter describes the design of a wavelet-based anonymity metric proposed in [12].

4.1 Wavelet Transformation

The wavelet transform is a mathematical tool that divides data or a function into different frequency components and then represents each component with a resolution matched to its scale. The most important advantage of the wavelet transform is that it converts the input signal into frequency and time localization. Time-frequency localization allows

studying the frequency contents locally in time. So the wavelet transform contains information about the frequency component at a particular point in time. The wavelet transform can be performed on a continuous-time signal as well as discrete-time signal. The wavelet transform for a continuous-time signal is as follows (as defined in [13])

$$(T^{wav} f)(a, b) = |a|^{-1/2} \int dt f(t) \psi \left(\frac{t-b}{a} \right) \quad (4.1)$$

and

$$(T_{m,n}^{wav} f) = a_0^{-m/2} \int dt f(t) \psi (a_0^{-m} t - nb_0) \quad (4.2)$$

In both cases it is assumed that ψ satisfies

$$\int dt \psi(t) = 0 \quad (4.3)$$

The second equation of wavelet transformation is obtained by restricting a, b to only a discrete set of values and $a = a_0^m, b = nb_0 a_0^m$, with m, n ranging over Z , and $a_0 > 1, b_0 > 0$ fixed. The following function $\psi^{a,b}$ is called the wavelet function:

$$\psi^{a,b}(s) = |a^{-1/2}| \psi \left(\frac{s-b}{a} \right) \quad (4.4)$$

In the case of the discrete wavelet transform, a and b can take only integers values.

Hence, $a = a_0^m$ and $a_0 > 1$. Different values of m correspond to wavelets of different widths. Narrow wavelets (high frequency) are translated by small steps in order to cover the whole time range, while wider wavelets (lower frequency) are translated by larger steps. The corresponding discrete wavelet form is as follows [13]:

$$\begin{aligned} \psi_{m,n}(x) &= a_0^{-m/2} \psi((a_0^{-m}(x - nb_0 a_0^m))) \\ &= a_0^{-m/2} \psi((a_0^{-m} x - nb_0)) \end{aligned} \quad (4.5)$$

4.1.1 Haar Wavelet

For some special choices of ψ and $a_0 = 2, b_0 = 1$, there exist ψ with good time-frequency localization properties [13], such that

$$\psi_{m,n}(x) = 2^{-m/2} \psi(2^{-m}x - n) \quad (4.6)$$

The Haar wavelet is also called the Daubechies 2 wavelet. The function $\psi_{m,n}(x)$ constitutes an orthonormal basis of $L^2(R)$. This function is called the Haar function [13],

$$\psi(x) = \begin{cases} 1, & 0 \leq x \leq \frac{1}{2} \\ -1, & \frac{1}{2} \leq x \leq 1 \\ 0, & \text{otherwise.} \end{cases}$$

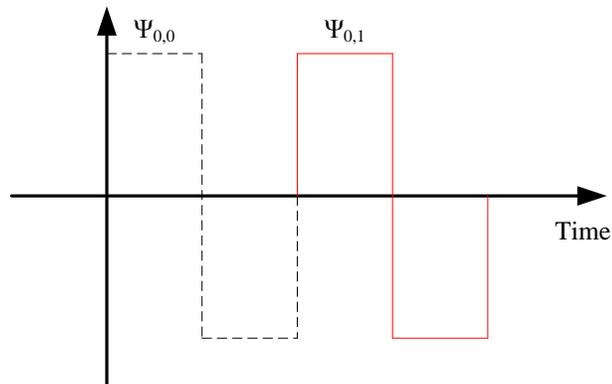


Fig.4.1 Haar Wavelet

4.2 Wavelet-based Anonymity Metric

This section describes a wavelet-based anonymity metric [12] that is used to determine the anonymity level of a particular network. The previous section explained the basics of

the wavelet transform. The timing constraints imposed by the low-latency requirement makes a network even more susceptible to timing attacks. Timing attacks exploit the timing information that is available by analyzing the incoming packets and the outgoing packets. With the real-time communication becoming faster, the end-to-end delay time reduces, which makes it even more difficult to eliminate the timing correlation between the flows. In order to measure anonymity, a generic metric is needed such that it is not specific to a particular type of network or configuration. The method proposed in [12] involves a novel metric that can quantitatively measure the practical effectiveness of all anonymous networks. By this it is possible to compare the anonymity offered by different anonymous networks. The anonymity metric in [12] uses the packet timing distortion between the incoming packets and the outgoing packets. It also measures anonymity in the presence of mixing, splitting, merging and packet dropping. The entire procedure discussed in this section is taken from [12].

4.2.1 Network Assumptions

An energy-based anonymity metric is used to determine the practical resilience of the network against timing attacks. The effectiveness of anonymous network as a variability of the packet timing distortion is converted to energy metric. The method in [12] uses the timing distortion and calculates the variability (randomness between the consecutive timing distortions) of the timing stamps. This is then used to determine the anonymity level in terms of energy. This energy metric is used as a measure of anonymity of that particular anonymous network.

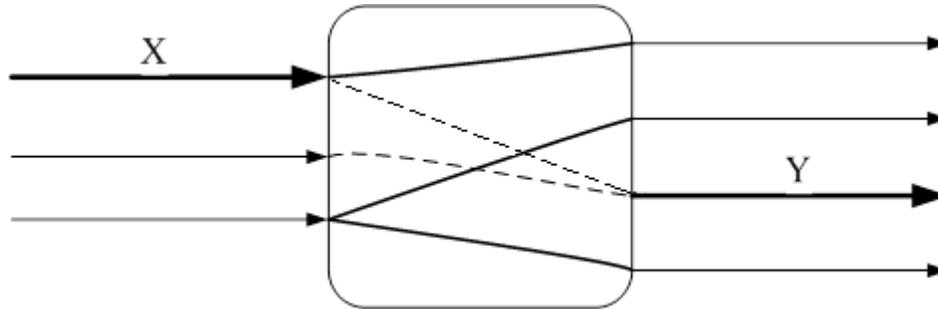


Fig. 4.2 Low Latency Anonymous Network

Figure 4.2 shows a low-latency anonymous network. The network is assumed to be a black box with no attackers inside the network. So only the incoming flows and outgoing flows are considered. The flows are also encrypted, which makes it difficult for the attacker to track the content of the flow. The incoming flow can go through various modifications like repacketization (combining several packets into one, splitting one packet into many), flow splitting, adding dummy packets and flow mixing. Let X be the incoming flow and Y be the corresponding outgoing flow. There may be many outgoing flows for a single incoming flow and many incoming flows for one outgoing flow due to the packet and flow modifications.

As said before, the network is assumed to be black box. The attacker can see the packets entering the network and the packets leaving the network. The attacker will not know anything about the flow that happens inside the network. The low-latency feature of the communication enables the attacker to gain some information from the timings of

the packets entering or leaving the network. Less the information about the packet timings implies more anonymity in the network and higher resilience to attacks.

4.2.2 Timing Distortion

Consider two flows that are to be analyzed. One is the incoming flow that enters the network. The other is the outgoing flow that leaves the network. Let X be the incoming flow to the anonymous network and Y be the corresponding outgoing flow. The method described in [12] is used to calculate the timing distortion from the timing information of the flows (incoming and outgoing).

Let T_f be used to represent the entire duration of the flow. Let n be the total number of packets in flow X and m be the total number of packets in flow Y . The numbers m and n need not be same. Let $P_0^x, P_1^x \dots \dots P_{n-1}^x$ represent the packets in flow X and $P_0^y, P_1^y \dots \dots P_{m-1}^y$ represent the packets in flow Y . Let $t(P_i^x), t(P_i^y)$ represents the time stamps of the i^{th} packet in each of flow X and flow Y respectively. The entire flow duration is divided into $\lceil T_f/T \rceil$ time intervals of equal lengths with $T > 0$. Let $S(i)$ denote the start time of interval i . The number of packets in interval i of flow f is denoted by $n(f, i)$. The term $\bar{t}(f, i)$ represents the mean time value of the packets in the interval i . When $n(f, i) = 0$, then $\bar{t}(f, i) = 0$.

Let,

$$x(f, i) = [\bar{t}(f, i) - \bar{t}(f, fpne(f, i))] \times n(f, i) \quad (4.7)$$

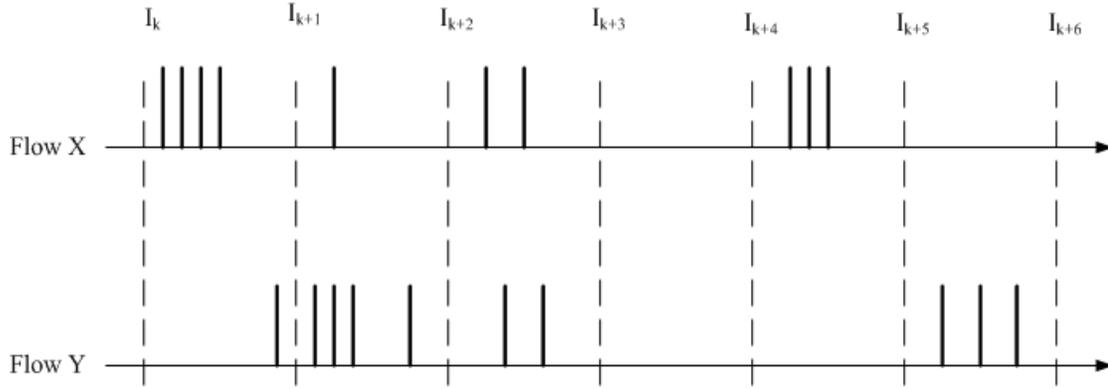


Fig. 4.3 Packet Timing Distortion

For interval i ($i > 0$) of flow f , $fpne(f, i)$ is defined as first previous non empty index that is before interval i . This will give the previous interval in which there were packets in them. For the first interval of flow f , $fpne(f, 0) = 0$. The aggregated time difference of interval i between flow f_1 and flow f_2 is defined as follows

$$d(f_1, f_2, i) = [x(f_1, i) - x(f_2, i)] \times S(i + 1) \quad (4.8)$$

Here $d(f_1, f_2, i)$ can be positive, negative or zero. An example of packet timing distortion is given in Figure 4.3. In this figure, flow X has fewer packets than flow Y in $(k + 1)$ th interval $I_{k+1} = [S(k + 1), S(k + 2))$ and $fpne(X, k + 1) = fpne(Y, k + 1) = k$. Hence, $d(X, Y, k + 1) < 0$. In the k th interval on the other hand, flow X has more packet arrivals than flow Y . Assuming for both the flows $fpne(X, k) = fpne(Y, k)$, then $d(X, Y, k + 1) < 0$. In the $(k + 3)$ th interval, flow X and flow Y does not have any packets in it and so $d(X, Y, k + 3) = 0$.

Therefore, the overall packet timing distortion vector that is between flow f_1 and flow f_2 is a given as follows [12],

$$D(f_1, f_2) = \langle d(f_1, f_2, 0), \dots \dots \dots, d\left(f_1, f_2, \left\lceil \frac{T_f}{T} \right\rceil - 1\right) \rangle \quad (4.9)$$

4.2.3 Multi Resolution Analysis

This section describes the analysis done on the variability of packet timing distortion between two flows using wavelet-based Multi Resolution Analysis. The method described in [12] uses a statistical estimator developed by Abry and Veitch [14]. The wavelet-based Multi Resolution analysis (MRA) [12] takes a sequence of data as input and transforms the sequence of data into a number of wavelet coefficients at different resolutions and at a low resolution approximation. For flow X and flow Y , interval size $T_0 > 0$ and the packet timing distortion vector is $\langle d(X, Y, 0), \dots, d\left(X, Y, \left\lceil \frac{T_f}{T_0} \right\rceil - 1\right) \rangle$. This vector is fed to wavelet based MRA to generate a series of vectors of different scales j :

$$D(X, Y, j) = \langle d_{j,0}, \dots \dots \dots, d_{j,n_j-1} \rangle$$

where, $n_j = \left\lceil \frac{T_f}{T_j} \right\rceil$,

$$T_j = 2^j T_0 \quad (j = 0, 1, \dots \dots \dots).$$

Also,

$$d_{j,k} = d_{j-1,2k} + d_{j-1,2k+1}, \quad j > 0$$

Let $C_{D(X,Y,j)}(p)$ be the p th ($p = 0, 1, \dots, N_j - 1$) wavelet detail coefficient at scale j for j^{th} vector $D(X, Y, j)$. The energy at time scale j is defined as the variance of the coefficients. The efficiency of the time averaging relies on the assumption that almost no correlation between the averaged quantities [15]. The averaged quantity here is $C_{D(X,Y,j)}(p)$. Therefore when $E \left(C_{D(X,Y,j)}(p) \right) = 0$ [15] is assumed, the energy at scale j is

$$e_j = \frac{\sum_{p=0}^{n_j-1} [C_{D(X,Y,j)}(p)]^2}{n_j} \quad (4.10)$$

Here, the wavelet based MRA assumes that $D(X, Y, j)$ is covariance stationary. This is meant for a given j , the mean of $D(X, Y, j)$ is a constant and the covariance between $d_{j,k}$ and $d_{j,k'}$ only depends on $|k - k'|$. The wavelet detail coefficient of the wavelet transform can be thought as an inner product of a high pass filter g (a vector of length 1) and the vector $\langle d_{j,2p}, \dots, d_{j,2p+l-1} \rangle$.

Using a Haar wavelet (Daubechies 2, $l = 2$) (section 4.1.1) [13] of scale j , whose high pass filter $g = \langle g_0, g_1 \rangle = \langle \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \rangle$:

$$\begin{aligned} C_{D(X,Y,j)}(p) &= \frac{1}{\sqrt{2^{j-1}}} g \cdot \hat{d}_{j-1} \\ &= \frac{1}{\sqrt{2^{j-1}}} (g_0 d_{j-1,2p} + g_1 d_{j-1,2p+1}) \\ &= \frac{1}{\sqrt{2^j}} (d_{j-1,2p} - d_{j-1,2p+1}) \end{aligned} \quad (4.11)$$

From the above equation, it is effectively a difference between an even-numbered element and an odd-numbered element of the $(j - 1)$ th scale vector. Let $\Delta d_{j-1,p} = d_{j-1,2p} - d_{j-1,2p+1}$, the energy e_j at scale j for the Haar wavelet becomes

$$e_j = 2^j \frac{\sum_{p=0}^{n_j} \Delta d_{j-1,p}^2}{n_j} \quad (4.12)$$

Here the term $\Delta d_{j-1,p}$ represents the correlation function and the expected value of this term is assumed to be 0 [15]. Since $E(\Delta d_{j-1,p}) = 0$, the energy e_j at scale j can be taken as the variance of the data variation $\Delta d_{j-1,p}$. Similarly, the Daubechies 6 wavelet transform [12] uses $\mathbf{g} = \langle g_0, g_1, g_2, g_3, g_4, g_5 \rangle$. The p th wavelet coefficient at scale j is

$$\begin{aligned} C_{D(X,Y,j)}(p) &= \frac{1}{\sqrt{2^{j-1}}} \mathbf{g} \cdot \hat{\mathbf{d}}_{j-1} \\ &= \frac{1}{\sqrt{2^{j-1}}} (\sum_{q=0}^5 g_q d_{j-1,2p+q}) \end{aligned} \quad (4.13)$$

where $\hat{\mathbf{d}}_{j-1} = \langle d_{j-1,2p}, d_{j-1,2p+1}, \dots, d_{j-1,2p+5} \rangle^T$. Since $\sum_{k=0}^5 g_k = 0$ and for all $k \neq k'$, $E(d_{j-1,k}) = E(d_{j-1,k'})$ because $E(\Delta d_{j-1,p}) = 0$ [15]. The energy at scale j for the D6 wavelet becomes

$$e_j = 2^{-j+1} \frac{\sum_{p=0}^{n_j} (\sum_{q=0}^5 g_q d_{j-1,2p+q})^2}{n_j} \quad (4.13)$$

4.2.4 Properties of the Energy-Based Metric

For flow X, Y and time interval size T_0 , the packet timing distortion is obtained by calculating $D(X, Y)$. Let $e_j(D(X, Y))$ denote the energy between flow X and flow Y at scale j . This energy-based metric has the following properties [12]:

1. If there is no distortion then the energy $e_j(D(X, X)) = 0$ for all j . This is because, $D(X, X) = \langle 0, \dots, 0 \rangle$ and hence, $C_{D(X, Y)}(p) = 0$ for all j and p .
2. For all flow X and flow Y , $e_j(D(X, Y))$ satisfies commutative property i.e. $e_j(D(X, Y)) = e_j(D(Y, X))$.
3. There is no change in energy by adding a constant to the distortion. This means that $e_j(D(X, Y)) = e_j(D(Y, X) + \hat{c})$, where $\hat{c} = \langle c, \dots, c \rangle$ be any vector of constant c of the same number of elements as that of $D(X, Y)$.
4. There will be a constant change in the plot if the distortion is multiplied by a constant. If each element of $D(X, Y)$ is multiplied with $a \neq 0$, then $e_j(aD(X, Y)) = a^2 e_j(D(Y, X))$. In other words, multiplying the distortion by a non-zero constant will move the plot up or down by a constant.

Chapter 5 Experimental Results

As discussed in the previous chapter, wavelet coefficients are calculated from the timing distortion obtained from the packet time stamps as follows:-

$$C_{D(X,Y,j)}(p) = \frac{1}{\sqrt{2^j}} (d_{j-1,2p} - d_{j-1,2p+1})$$

and,

$$C_{D(X,Y,j)}(p) = \frac{1}{\sqrt{2^{j-1}}} (\sum_{q=0}^5 g_q d_{j-1,2p+q})$$

$$e_j = \frac{\sum_{p=0}^{n_j-1} [C_{D(X,Y,j)}(p)]^2}{n_j}$$

This chapter discusses the simulation results that were obtained above using the wavelet-based model. The wavelet-based energy plot shows the logarithm of energy $\log_2(e_j)$ at all time scales (X-axis), which shows the variability of the input sequence of data at different time scales. The more variable the output is the higher the energy will be. The resolution index is made as Y-axis. The resolution index indicates the reference time scale. The minimum reference time scale is taken as 1 second. As the resolution index increase, the time scale is multiplied by a factor of 2^j and j denotes the resolution index. In the Multi-resolution Analysis method (MRA), the resolution index is used to specify the period or duration of packet flow under consideration.

5.1 Energy Plots for Poisson Process

In the simulation model, packet arrivals follow a Poisson arrival process. The inter-arrival time periods, i.e., the difference in time between one packet and the next packet follows an exponential distribution for Poisson arrivals. The Poisson process is discussed in Section 3.2.1. The first plot analyzes the energy for different Poisson arrival rates. It consists of energy plots for three different rates or three different inter-arrival times. The average inter-arrival time used for simulation purposes are 0.9, 1.5 and 2.1 seconds. The Figure 5.1 shows three different colored lines corresponding to three different arrival rates. It can be seen that as the arrival rate decreases, the energy increases. This means the anonymity level is higher for higher inter-arrival times. The reason for this is, as the packet arrival rate increases, the inter-arrival time between the packets is less. In the case of small inter-arrival times, it is difficult to anonymize the packets and hence, the energy decreases. In the case of higher inter-arrival times, the network has more time to anonymize the packets and hence, higher energy levels. From Figure 5.1, the inter-arrival time of 2.1 seconds results in the maximum energy level and the energy for the inter-arrival time of 0.9 seconds is the lowest.

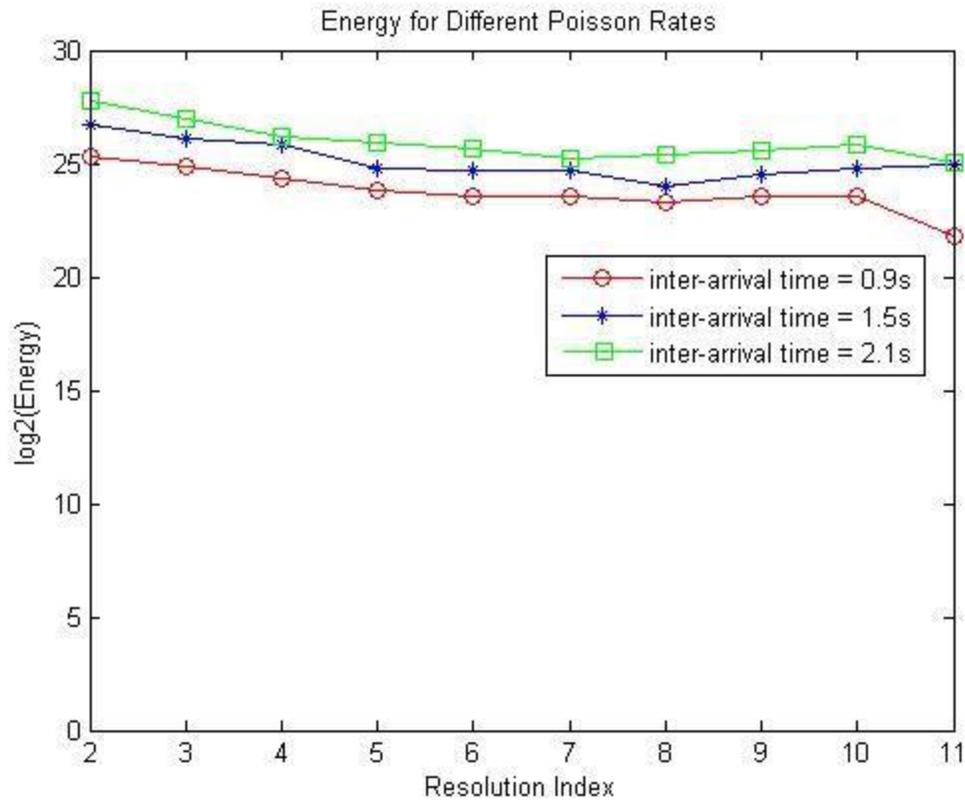


Fig. 5.1 Energy plot for 3 different Poisson Rates

Figure 5.2, 5.3 and 5.4 shows the energy plot for a fixed Poisson rate but different end-to-end delays. The end-to-end delay is generated using three models: The first is a constant delay, the second model characterizes the delay as a uniform random variable, and the third model characterizes the delay as a random process (autoregressive process). The constant delay is assumed to be 0.09 seconds. This means that the time taken for a packet to reach the destination user from its starting point takes 0.09 seconds. In the case of the delay as uniform random variable, it is generated with a mean value of 0.09 seconds, equal to constant delay. An autoregressive process is used to generate a random

process of delay values and the following equation is used to generate the end-to-end delay.

$$t_{delay}(i) = 0.8 * t_{delay}(i - 1) + 0.2 * mean + error \quad (5.1)$$

Here, *mean* represents the mean value of the end-to-end delay that was assumed for the constant delay or the mean of the uniform random variable (i.e., 0.09 seconds). The error term is generated using a normal random variable with variance 0.02 and zero mean. The term $t_{delay}(i)$ is used to calculate the i^{th} time delay from the previous time delay ($i - 1$). For $t_{delay}(1)$, it is assumed to take a random value between 0 and 0.09. Figure 5.2 shows the energy for 0.9 second inter-arrival time. The energy using the constant delay has the least value when compared to the other delay models. The energy using uniform random variable delay and autoregressive process are close to each other. This means the randomness of the packet time stamps, created by these delay models are very close. Therefore, the energy calculated is close to each other and hence, they have similar anonymity levels.

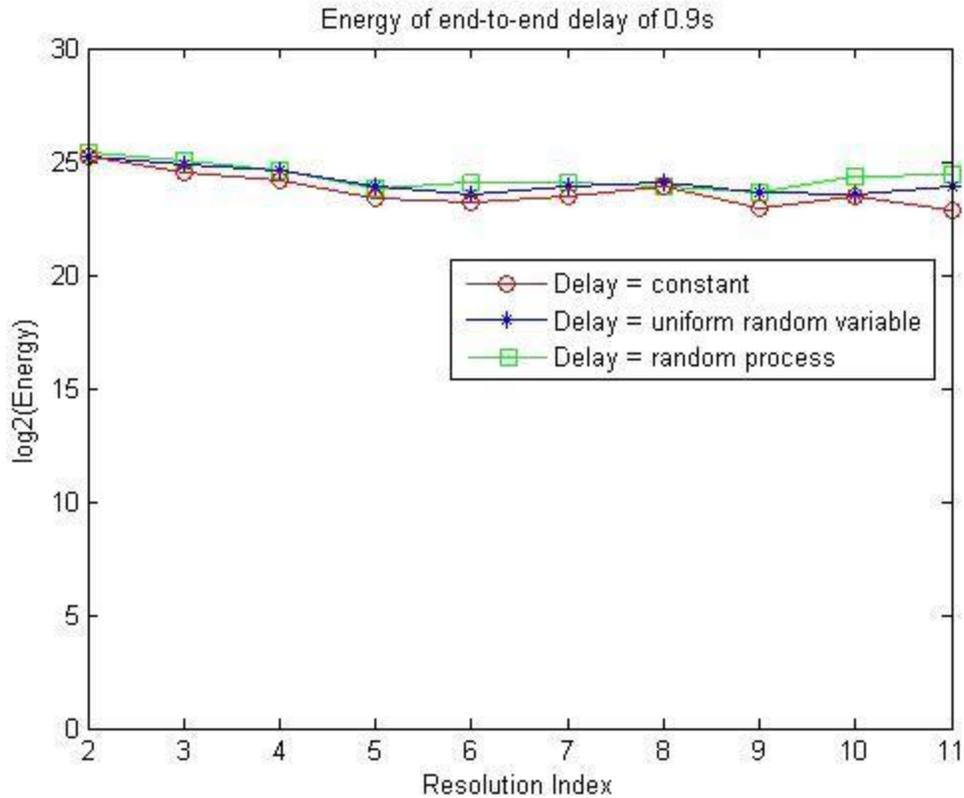


Fig. 5.2 Energy plot for different delays for inter-arrival time of 0.9 seconds

The next two graphs show the energy for an inter-arrival time of 1.5 seconds and 2.1 seconds. In each of the cases, constant delay between the end users has slightly less energy than the uniform distribution or an autoregressive process model for the end-to-end delay. Hence, the anonymity levels achieved for constant delay is slightly less when compared to uniform random variable delay and random process based delay.

The resolution index refers to the size of the period that is being analyzed. As discussed, the simplest resolution index is 1 and the size of period corresponding to this resolution index is 1 second. For each increment in the index, the value of the period is doubled. As for the resolution index 11, the size of one period corresponds to 1024

seconds. In Figures 5.2, 5.3 and 5.4, the energy for the uniform random variable and random process delay models deviates more during the measurement for index 11. This deviation creates a notable difference between the energy for constant delay and other inter-arrival models. As the resolution index increases, the reference time period increases. The time period for index 11 is the highest in resolution index used in this experiment. This requires a large set of values to calculate the energy. Since this is a large set of data, the randomness for the whole set of data appears to be less for the Poisson model compared to other inter-arrival models. Hence, there is a considerable difference in the energy levels at index 11.

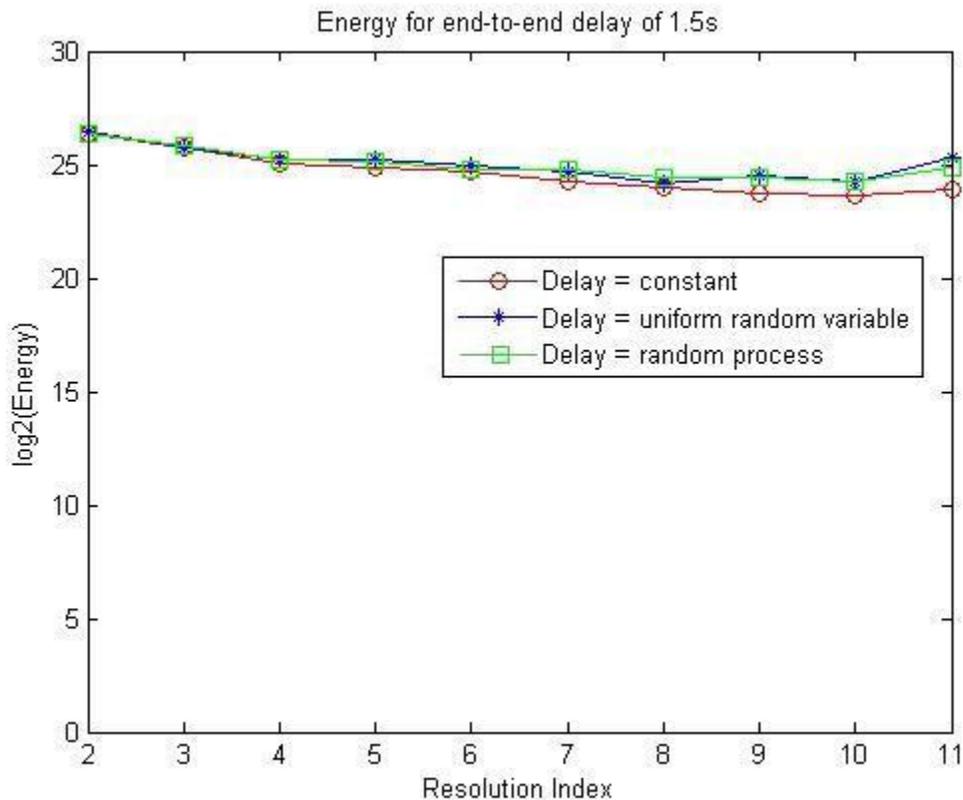


Fig. 5.3 Energy plot for different delays for inter-arrival time of 1.5 seconds

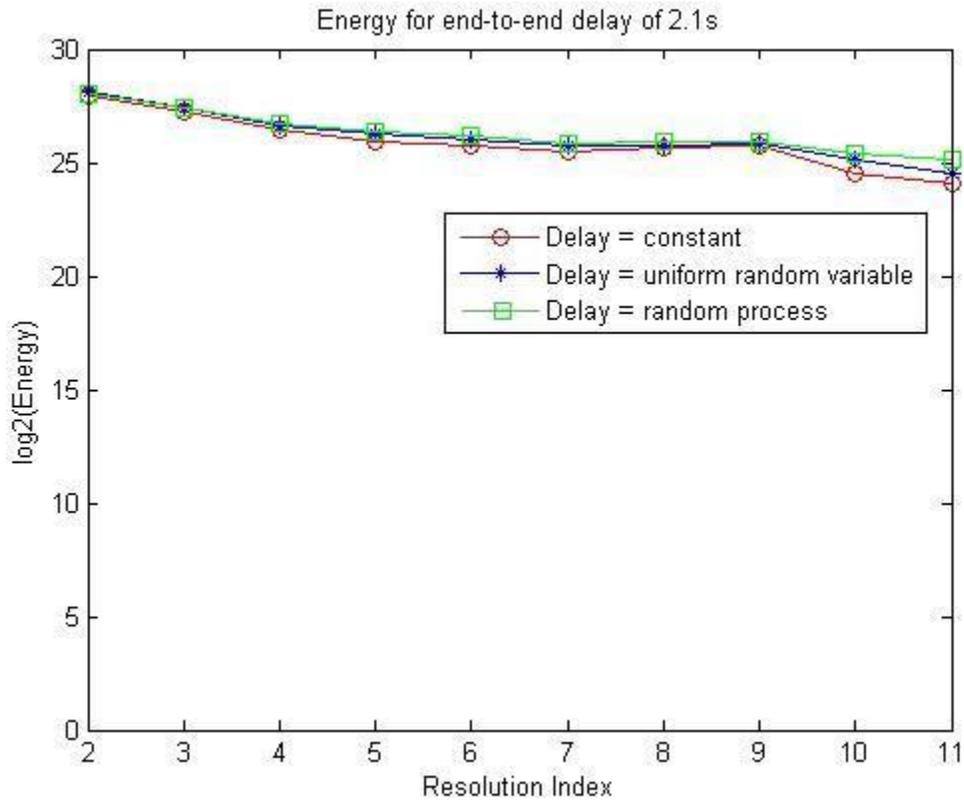


Fig. 5.4 Energy plot for different delays for inter-arrival time of 2.1 seconds

Figures 5.2, 5.3 and 5.4 show the energy for inter-arrival times of 0.9 seconds, 1.5 seconds and 2.1 seconds. For each of these inter-arrival times, the end-to-end delay is characterized as a constant, a random variable and an autoregressive process. The results show that for constant delay the energy is less than that of the delay characterized by the uniform random variable or the autoregressive process. This is because, in the anonymity level calculation, the distortion in the packet timings generated by the other random delay models is higher. Hence, the uniform random variable delay and random process delay models increase the anonymity level of the network relative to the constant delay model.

5.2 Energy Plots for Different Packet Arrival Models

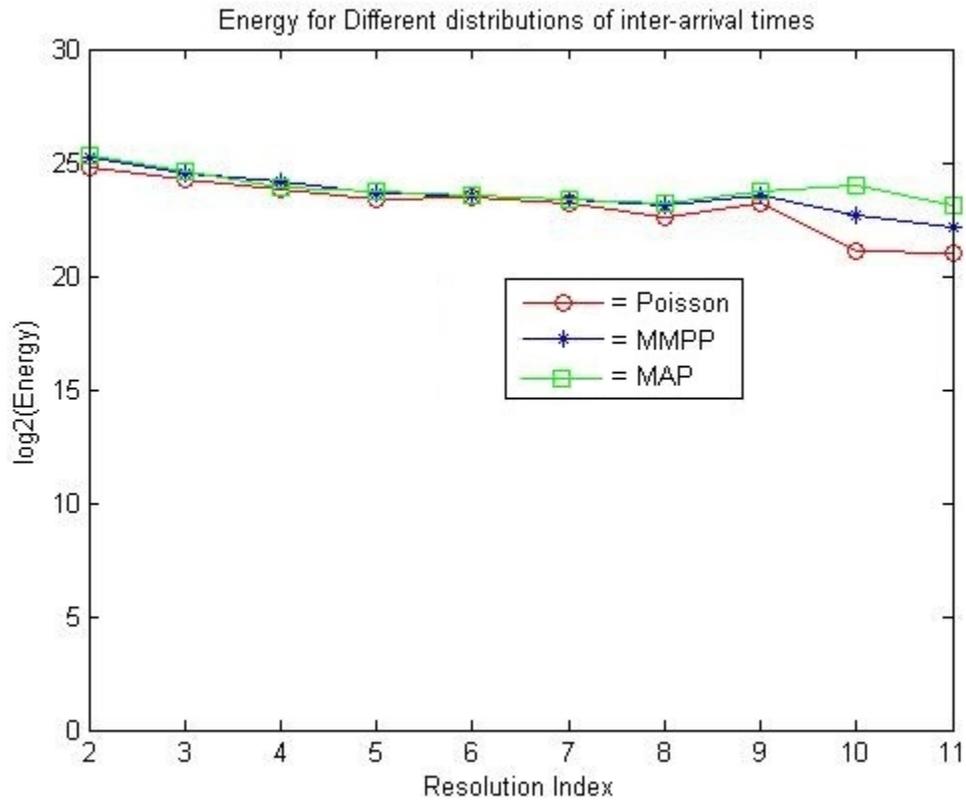


Fig. 5.5 Energy plot for different traffic models

Figure 5.5 shows the energy plot for different traffic models, namely the Poisson process, the Markov Modulated Poisson process (MMPP) and the Markovian Arrival Process (MAP).

5.2.1 Theoretical Calculation of Arrival Rate for MMPP AND MAP

In the case of MMPP, the Q (generator) matrix of the underlying chain is as follows:-

$$Q = \begin{bmatrix} -1.2 & 0.6 & 0.6 \\ 0.6 & 1.2 & 0.6 \\ 0.6 & 0.6 & -1.2 \end{bmatrix}$$

and,

$$\bar{\lambda} = (\lambda_1, \lambda_2, \lambda_3)$$

$$\bar{\lambda} = (0.9, 1.1, 1.3)$$

The theoretical calculation of mean arrival rate for the above rates is calculated to compare with the simulated value. Let π be the stationary probability vector of the Markov process.

$$\pi Q = 0 \text{ and } \pi E = \underline{1}. \quad (5.2)$$

Hence,
$$\pi(Q + E) = \underline{1}. \quad (5.3)$$

Here $\underline{1}$ is a row matrix of 1's and E is square matrix of all 1's.

$$\pi = \underline{1} * (Q + E)^{-1} \quad (5.4)$$

Solving this equation give the value for π .

$$\pi = |0.3333 \quad 0.3333 \quad 0.3333|$$

After applying the rates and solving π , the stationary arrival rate is

$$\lambda_{mean} = \langle \underline{\pi}, \underline{\lambda} \rangle$$

$$\lambda_{mean} = 1.1$$

In the case of MAP model, the generator for the bivariate process is as follows:-

$$G = \begin{bmatrix} D_0 & D_1 & 0 & \dots & \dots & \dots \\ 0 & D_0 & D_1 & 0 & \dots & \dots \\ 0 & 0 & D_0 & D_1 & 0 & \dots \\ 0 & 0 & \dots & D_0 & D_1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \ddots \end{bmatrix}$$

and,

$$D_0 = \begin{bmatrix} -2.3 & 0.5 & 0.5 \\ 0.6 & -2.1 & 0.6 \\ 0.6 & 0.6 & -2.1 \end{bmatrix},$$

$$D_1 = \begin{bmatrix} 0.3 & 0.45 & 0.35 \\ 0.3 & 0.45 & 0.35 \\ 0.3 & 0.45 & 0.35 \end{bmatrix}.$$

The theoretical value for the arrival rate in the case of MAP [16] is also calculated so as to compare with the simulated values. Let π be the stationary probability vector of the underlying Markov process with generator matrix $Q = D_0 + D_1$ and π satisfies the following equations.

$$\pi Q = 0 \text{ and } \pi E = \underline{1}. \quad (5.5)$$

E is a square matrix of all 1's.

$$\pi(Q + E) = \underline{1}$$

$$\pi = \underline{1} * (Q + E)^{-1} \quad (5.6)$$

$$\pi = [0.2857 \quad 0.3764 \quad 0.3379]$$

The above equation calculates the stationary probability vector. The arrival rate is

$$\lambda_{mean} = \pi D_1 e$$

After substituting the values, the mean arrival rate is $\lambda_{mean} = 1.1$ packet/second.

5.3 Energy Plots for MMPP and MAP Traffic Models

From the plot, it can be seen that the energy for MMPP and MAP inter-arrival times have higher energy than that of Poisson inter-arrival time. With the Markovian Arrival Process and Markov Modulated Poisson Process models, the inter-arrival times create more timing distortion than that of Poisson inter-arrival time. The anonymity level of a network would be higher if packet arrivals are modelled with MAP or MMPP.

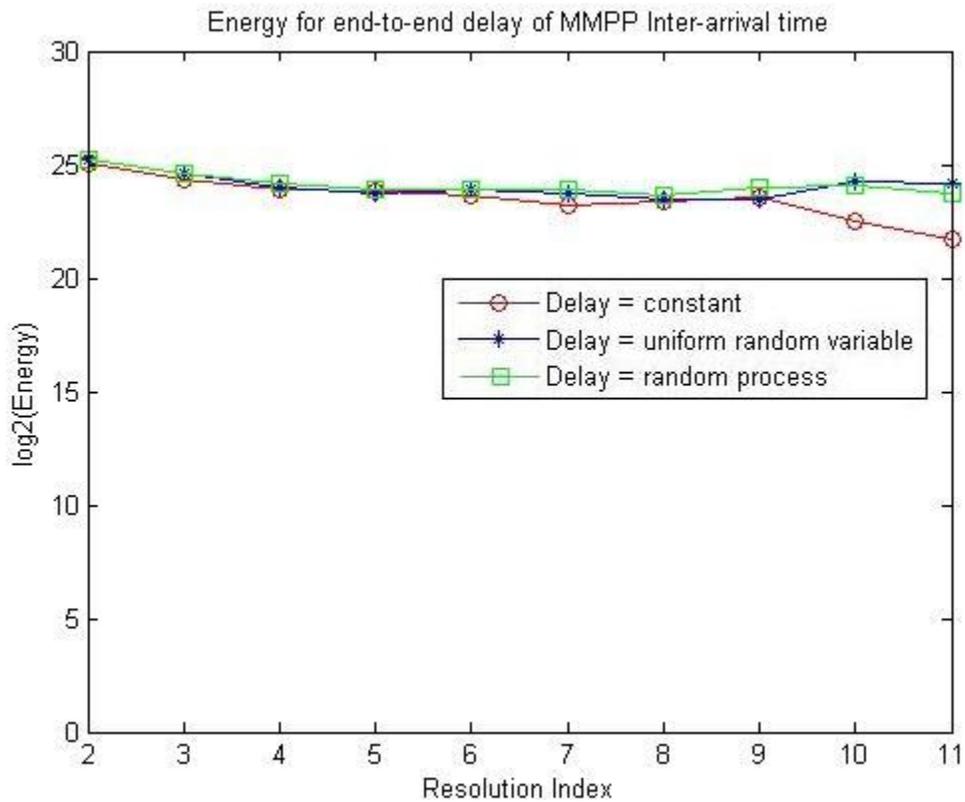


Fig. 5.6 Energy plot for different delays for MMPP

Figure 5.6 shows the energy for the MMPP traffic model with different models of end-to-end delay. In this plot, the energy levels for the different delays are very close to each other unlike the case with Poisson process model. The inter-arrival times using MMPP model itself creates much of the timing distortion when compared to end-to-end timing delay. Eventhough they are close to each other, the curves for the autoregressive model and the uniform random variable model were slightly higher than that of the constant delay model.

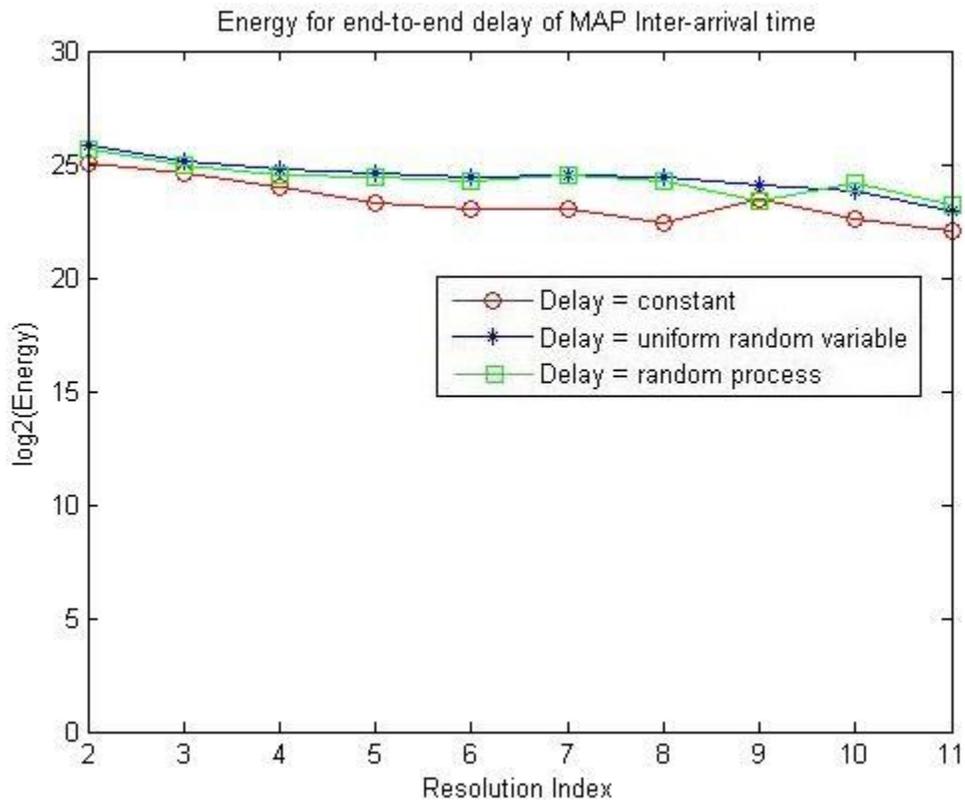


Fig. 5.7 Energy plot for different delays for MAP.

Figure 5.7 shows the energy for MAP based inter-arrival time with different models of end to end delay. In this plot, the energy levels for the different delays are similar to that of MMPP. The energy level for the three delay models are close to each other, which shows that the anonymity levels are similar for different delay models when different traffic models are used. This is because there is not much difference in the packet timing distortions can be achieved by changing the end-to-end delay when the traffic pattern is more complex.

5.4 Conclusion

Using packet time stamps, the anonymity level of a network is measured empirically via a quantitative anonymity metric. In the first part of the experiment, different Poisson rates were used to determine the energy plot changes. Understanding that the energy plot changes with change in the packet arrival rates is of great importance. From Figure 5.1, it can be seen that as the inter-arrival times increases (decreasing Poisson rates), the energy increases. An increase in energy implies that the measured anonymity is higher for higher inter-arrival time between packets.

For the second part of the experiment, different delay models are used for Poisson, MMPP and MAP arrival process. The different delay models used are: constant delay, uniform random delay and autoregressive random process delay. The energy for a particular traffic model is lower for constant delay than the other two delay models. The last part of the experiment compared the energy achieved for different traffic models.

From Figure 5.5, it can be seen that the energy for MMPP and MAP is higher than that of Poisson arrival rate.

The energy measure provides a quantitative characterization of the anonymity level in a network. The anonymity is higher if the packet timing distortion is high, i.e., if the time stamps are more random. The randomness created by the MAP and MMPP is more than that of the Poisson model. So the energy for MAP and MMPP were higher which means that the anonymity is higher if the packet traffic follow MAP and MMPP models.

Similarly, with the end-to-end delay, the uniform random variable model and random process model caused the timing distortion to be much more random, resulting in a higher anonymity level compared to that of constant delay. There is only a small difference between the energy attained for MAP and MMPP models because the timing distortions generated by these two models were substantially different.

Chapter 6 Conclusion

A wavelet-based energy metric proposed in [12], provides a means of calculating the anonymity level of a network. The method uses packet time stamps to calculate the energy value. Many methods have been proposed to calculate the anonymity of a network. The method of [12], which is the focus of this thesis, is an effective method for quantifying anonymity. Some of the previous methods that were based on information theory had some drawbacks. In calculating the information-theoretic anonymity metric, the attacker has to use probabilistic assumptions. The probabilistic assumptions are user-specific and for the same network the calculated anonymity value may be different for each attacker. One other drawback is that if the network setup changes (addition or removal of a user, router, etc.), all the information collected thus far will not be of any use. The wavelet-based energy metric is independent of the network setting and is generic in nature. Hence, it is very effective for calculating anonymity even for a newly designed anonymous network.

In this thesis, we analyzed different packet arrival models and calculated the wavelet-based anonymity metric for each one of them. For each of these packet arrival models, the end-to-end delay is characterized using different probability models. The time stamps of the packets, created using complex distributions, were more random and variable. The energy calculated using this method was higher for complex packet arrival distributions.

Hence, the anonymity calculated for a network, keeping the setup as same, will be higher if the packet arrivals follow more complex traffic models and the time stamps of packets are more random.

Bibliography

Bibliography

- [1] M. G. Reed, P. F. Syverson and D. M. Goldschlag. Anonymous Connections and Onion Routing. *IEEE JSAC Copyright and Privacy Protection*, 16(4):482-494, 1998.
- [2] O. Berthold, H. Federrath and S. Kpösell. Web MIXes: A system for Anonymous and Unobservable Internet Access. In *Proceedings of Designing Privacy Enhancing Technologies: Workshop on Design Issues in Anonymity and Unobservability*, pages 115-129, July 2000.
- [3] B. N. Levine, M. K. Reiter, C. Wang and M. K. Wright. Timing Analysis in Low-Latency Mix-Based Systems. In A. Juels, editor, *Proceedings of Financial Cryptography (FC '04)*, pages 251-265. Springer-Verlag, LNCS 3110, February 2004.
- [4] A. Srgantov and G. Danezis. Towards an information theoretic metric for Anonymity. In *Proceedings of the 2nd International Workshop of Privacy Enhancing Technologies workshop (PET 2002)*, pages 41-53, San Francisco, CA, USA, 2002. Springer-Verlag.
- [5] C. Diaz, S. Seys, J. Claessens and B. Preneel. Towards measuring anonymity. In *Proceedings of the 2nd International Workshop of Privacy Enhancing Technologies workshop (PET 2002)*, pages 54-68, San Francisco, CA, USA, 2002. Springer-Verlag.
- [6] B.N. Levine, M.K. Reiter, C. Wang, and M.K. Wright. Timing Attacks in Low-Latency Mix Systems. In A. Juels, editor, *Proceedings of Financial Cryptography (FC 04')*, pages 251-265. Springer-Verlag, LCNS 3110, February 2004.
- [7] S.J. Murdoch and G. Danezis. Low-cost Traffic Analysis of TOR. In *Proceedings of 2005 IEEE Symposium on Security and Privacy (S & P 2005)*, pages 183-195, May 2005.
- [8] R. Dingledine, N. Mathewson, and P. Syverson. TOR: The Second-Generation Onion Routing. In *Proceedings of the 13th USENIX Security Symposium*, pages 303-320, San Diego, CA, USA, August 2004. USENIX.
- [9] E. Cinlar, *Introduction to Stochastic Processes*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1975.

- [10] W. Fischer, K.M. Hellstern. The Markov-Modulated Poisson Process (MMPP) cookbook. In *Proceedings of Performance Evaluation, volume 18, issue 2*, pages 149-171. March 1991.
- [11] D.M. Lucantoni, New Results on the Single Server Queue with a Batch Markovian Arrival Process. *Stochastic Models, volume 7, No. 1*, pages 1-46, 1991.
- [12] J. Jin and Xinyuan Wang, On the Effectiveness of Low latency Anonymous Network in the Presence of Timing Attacks. In *Proceedings of the 39th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DNS 2009)*, July 2009.
- [13] I. Daubechies, *Ten Lectures on Wavelets*. Philadelphia, PA: SIAM, 1992.
- [14] P. Abry and D. Weitch. Wavelet Analysis of Long-Range Dependent Traffic. *IEEE Transactions on Information Theory*, 44(1): 2-15, January 1998.
- [15] P. Abry, P. Gonçalvès, and P. Flandrin, “Wavelet-based spectral analysis of $1/f$ processes,” in *Proceedings IEEE-ICASSP '93, 1993*, pp. III.237-III.240.
- [16] David M. Lucantoni, “The BMAP/G/1 Queue: A Tutorial”, in *Proceedings Performance Evaluation of Computer Communication Systems, Joint Tutorial Papers of Performance '93 and sigmetrics '93*, pages 330-358. Springer, Berlin, 1993.

Curriculum Vitae

Abinash Vasudevan received Bachelor of Engineering in Electronics and Communication Engineering in 2009 from Saveetha Engineering College. He then pursued Master of Science in Electrical Engineering in the Department of Electrical and Computer Engineering at George Mason University. He worked as a graduate teaching assistant for the Department of Electrical and Computer Engineering.