

On Flow Correlation Attacks and Countermeasures in Mix Networks^{*}

Ye Zhu*, Xinwen Fu, Bryan Graham*, Riccardo Bettati, and Wei Zhao

Department of Computer Science, Texas A&M University,
College Station TX 77843-3112, USA
{zhuye, bgraham}@tamu.edu,
{xinwenfu, bettati, zhao}@cs.tamu.edu

Abstract. In this paper, we address issues related to flow correlation attacks and the corresponding countermeasures in mix networks. Mixes have been used in many anonymous communication systems and are supposed to provide countermeasures that can defeat various traffic analysis attacks. In this paper, we focus on a particular class of traffic analysis attack, *flow correlation attacks*, by which an adversary attempts to analyze the network traffic and correlate the traffic of a flow over an input link at a mix with that over an output link of the same mix. Two classes of correlation methods are considered, namely *time-domain* methods and *frequency-domain* methods. Based on our threat model and known strategies in existing mix networks, we perform extensive experiments to analyze the performance of mixes. We find that a mix with any known batching strategy may fail against flow correlation attacks in the sense that for a given flow over an input link, the adversary can correctly determine which output link is used by the same flow. We also investigated methods that can effectively counter the flow correlation attack and other timing attacks. The empirical results provided in this paper give an indication to designers of Mix networks about appropriate configurations and alternative mechanisms to be used to counter flow correlation attacks.

1 Introduction

This paper studies flow correlation attacks and the corresponding countermeasures in mix networks. With the rapid growth and public acceptance of the Internet as a means of communication and information dissemination, concerns about privacy and security on the Internet have grown. Although it can potentially be used for malicious purposes, *Anonymity* is legitimate in many scenarios such as anonymous web browsing, E-Voting,

^{*} This work was supported in part by the National Science Foundation under Contracts 0081761 and 0324988, by the Defense Advanced Research Projects Agency under Contract F30602-99-1-0531, and by Texas A&M University under its Telecommunication and Information Task Force Program. Any opinions, findings, and conclusions or recommendations in this material, either expressed or implied, are those of the authors and do not necessarily reflect the views of the sponsors listed above.

E-Banking, E-Commerce, and E-Auctions. In each of these scenarios, encryption alone cannot achieve the anonymity required by participants [1, 2].

Since Chaum [3] proposed the mix network, researchers have developed various anonymity systems for different applications. Although a significant amount of effort has been put forth in researching anonymous communications, there has not been much systematic study of the performance of mix networks in terms of anonymity degree provided and quality-of-services maintained. This paper focuses on the quantitative evaluation of mix performance. We are particularly interested in flow-based communication, which is widely used in voice over IP, web browsing, FTP, etc. These applications may have anonymity requirements, and the mixes are supposed to provide countermeasures that can defeat traffic analysis attacks.

We focus our analysis on a particular type of attack, which we call a *flow correlation attack*. In this type of attack, an adversary analyzes the network traffic with the intention of identifying which of several output ports a flow at an input port of a mix is taking. Obviously, flow correlation helps the adversary identify the path of a flow and consequently reveal other mission critical information related to the flow (e.g., sender and receiver). Our major contributions are summarized as follows:

1. We formally model the behavior of an adversary who launches flow correlation attacks. In order to successfully identify the output port of an incoming flow, the flow correlation attack must accurately measure the similarity of traffic flows into and out of a mix. Two classes of correlation methods are considered, namely *time-domain* methods and *frequency-domain* methods. In the time domain, *mutual information* is used to measure the traffic similarity. In the frequency domain, a matched filter based on the *Fourier spectrum* and the *Wavelet spectrum* is utilized.
2. We measure the effectiveness of a number of popular mix strategies in countering flow correlation attacks. Mixes with any tested batching strategy may fail under flow-correlation attacks in the sense that, for a given flow over an input link, the adversary can effectively detect which output link is used by the same flow. We use *Detection rate* as the measure of success for the attack, where Detection rate is defined as the probability that the adversary correctly correlates flows into and out of a mix. We will show that, given a sufficient amount of data, known mix strategies fail, that is, the attack achieves close to 100% detection rate. This remains true, even in batching strategies that sacrifice QoS concerns (such as a significant TCP goodput reduction) in favor of security.
3. While many mix strategies rely on other mechanisms in addition to batching alone, it is important to understand the vulnerability of batching. In our experiments, we illustrate the dependency between attack effectiveness for various batching strategies and the amount of data at hand for the attacks. These results should guide mix designers in the educated choice of strategy parameters, such as for striping or for path rerouting.

To counter flow correlation attacks, we investigate countermeasures based on theoretical analysis. We purposely synchronize the sending time of packets along a set of output links. This approach is more efficient than similar methods.

The remainder of this paper is organized as follows: Section 2 reviews the related work. Section 3 outlines our Mix network model, threat model, and a formal definition

of the problem. Batching strategies used by existing mix networks are also discussed in this section. Section 4 introduces traffic analysis methodologies that may be deployed by an adversary. We consider both time-domain and frequency-domain traffic analysis methods. In Section 5 we evaluate the performance of mix networks in terms of detection rate and FTP goodput. Serious failure of mix networks in terms of providing flow anonymity is observed from the data we collect. In Section 6, we present an effective and efficient method that can provide a guaranteed detection rate with high FTP goodput. We conclude this paper and discuss the future work in Section 7.

2 Related Work

Chaum [3] pioneered the idea of anonymity in 1981. Since then, researchers have applied the idea to different applications such as message-based email and flow-based low-latency communications, and they have invented new defense techniques as more attacks have been proposed.

For anonymous email applications, Chaum [3] proposed to use relay servers, i.e. *mixes*, rerouting messages, which are encrypted by public keys of mixes. An encrypted message is analogous to an onion constructed by a sender, who sends the onion to the first mix. Using its private key, the first mix peels off the first layer, which is encrypted using the public key of the first mix. Inside the first layer is the second mix's address and the rest of the onion, which is encrypted with the second mix's public key. After getting the second mix's address, the first mix sends the peeled onion. This process proceeds in this recursive way. The core part of the onion is the receiver's address and the real message to be sent to the receiver by the last mix. Chaum also proposed return address and digital pseudonyms for users to communicate with each other anonymously.

Helsingius [4] implemented the first Internet anonymous *remailer*, which is a single application proxy that just replaces the original email's source address with the remailer's address. It has no reply function and is subject to all the attacks mentioned below. Eric Hughes and Hal Finney [5] built the *cypherpunk remailer*, a real distributed mix network with reply functions that uses PGP to encrypt and decrypt messages. The system is subject to a global passive attack and replay attack to its reply mechanism. Gülcü and Tsudik [6] developed a relatively full-fledged anonymous email system, *Babel*. Their reply technique does not need the sender to remember the secret seed to decrypt the reply message, but it is subject to replay attack. They studied the threat from the trickle attack, a powerful active attack. Another defect of Babel is that a mix itself can differentiate the forwarding and replying messages. Cottrell [7] developed *Mixmaster* which counters a global passive attack by using message padding and also counters trickle and flood attacks [6, 8] by using a pool batching strategy. Mixmaster does not have a reply function. Danezis, Dingle and Mathewson [9] developed *Mixminion*. Although Mixminion still has many problems, its design considers a relatively complete set of attacks that researchers have found [8, 10, 11, 12, 13, 14]. The authors suggest a list of research topics for future study.

Low-latency anonymous communication can be further divided into systems using core mix networks and peer-to-peer networks. In a system using a core mix network, users connect to a pool of mixes, which provides anonymous communication, and users

select a forwarding path through this core network to the receiver. *Onion routing* [15] and *Freedom* [16] belong to this category. In a system using a peer-to-peer network, every node in the network is a mix, but it can also be a sender and receiver. Obviously, a peer-to-peer mix network can be very large and may provide better anonymity in the case when many participants use the anonymity service and enough traffic is generated around the network. *Crowds* [17], *Tarzan* [18] and *P⁵* [19] belong to this category.

This paper is interested in the study of passive traffic analysis attacks against low-latency anonymous communication systems. Sun *et al.* [2] gave a quantitative analysis for identifying a web page even if encryption and anonymizing proxies are used. They took advantage of the fact that a number of HTTP features such as the number and size of objects can be used as signatures to identify web pages with some accuracy. Unless the anonymizer addresses this, these signatures are visible to the adversary. Serjantov and Sewell [20] analyzed the possibility of a lone flow along an input link of a mix. If the rate of this lone input flow is roughly equal to the rate of a flow out of the mix, this pair of input flow and outflow flow are correlated. They also briefly discussed some of the possible traffic features used to trace a flow. The attacks we will present later in this paper are very effective even when a large amount of noise exists. Other analyses focus on the anonymity degradation when some mixes are compromised, e.g. [17]. We understand that attacks used against message-based email mix networks can also threaten low-latency flow-based mix networks; however, we feel that traffic analysis attacks are also a serious problem for low-latency mix networks because of its QoS requirements. Our reasoning will be explained in detail in the following sections of this paper.

3 Models

3.1 Mix and Mix Network

A mix is a relay device for anonymous communication. Figure 1 shows the communication between users using one mix. A single mix can achieve a certain level of communication anonymity: The sender of a message attaches the receiver address to a packet and encrypts it using the mix's public key. Upon receiving a packet, a mix decodes the packet. Different from an ordinary router, a mix usually will not relay the received packet immediately. Rather, it collects several packets and then sends them out in a *batch*. The order of packets may be altered as well. Techniques such as batching and reordering are

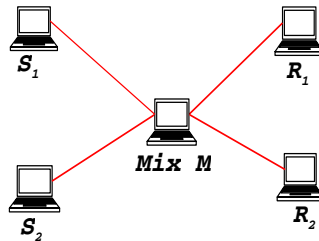


Fig. 1. A Single Mix

considered necessary techniques for mixes to prevent timing-based attacks. The main objective of this paper is to analyze the effectiveness of mixes against a special class of timing-based attacks.

A mix network consists of multiple mixes that are inter-connected by a network. A mix network may provide enhanced anonymity, as payload packets may go through multiple mixes. Even in such a mix network, it is important that each individual mix provides sufficient security and QoS so that the end-to-end performance can be guaranteed. Thus, our analysis on a single mix provides a foundation for analyzing the end-to-end performance of mix networks. We discuss in detail how to extend our work to larger and complicated mix networks in [21]. In fact, if we view a mix network (for example Onion routing [15]) as one *super mix*, the analytical techniques in this paper can be directly applied.

3.2 Batching Strategies for a Mix

Batching strategies are designed to prevent not only simple timing analysis attacks but also powerful trickle attacks, flood attacks, and many other forms of attacks ([9, 8]). Serjantov [8] summarizes seven batching strategies that have been proposed. We will evaluate each kind of these strategies. Our results show that these strategies may not work under certain timing analysis attacks. These seven batching strategies are listed in Table 1, in which batching strategies from S_1 to S_4 are denoted as *simple mix*, while batching strategies from S_5 to S_7 are denoted as *pool mix*.

Table 1. Batching Strategies

Glossary			
n	queue size		
m	threshold to control the packet sending		
t	timer's period if a timer is used		
f	the minimum number of packets left in the pool for pool Mixes		
p	a fraction only used in Timed Dynamic-Pool Mix		

Algorithms			
Strategy Index	Name	Adjustable Parameters	Algorithm
S_0	Simple Proxy	<i>none</i>	no batching or reordering
S_1	Threshold Mix	$\langle m \rangle$	if $n = m$, send n packets
S_2	Timed Mix	$\langle t \rangle$	if timer times out, send n packets
S_3	Threshold Or Timed Mix	$\langle m, t \rangle$	if timer times out, send n packets; elseif $n = m$ {send n packets; reset the timer}
S_4	Threshold and Timed Mix	$\langle m, t \rangle$	if (timer times out) and $(n \geq m)$, send n packets
S_5	Threshold Pool Mix	$\langle m, f \rangle$	if $n = m + f$, send m randomly chosen packets
S_6	Timed Pool Mix	$\langle t, f \rangle$	if (timer times out) and $(n > f)$, send $n - f$ randomly chosen packets
S_7	Timed Dynamic-Pool Mix	$\langle m, t, f, p \rangle$	if (timer times out) and $(n \geq m + f)$, send $\max(1, \lfloor p(n - f) \rfloor)$ randomly chosen packets

From Table 1, we can see that the sending of a batch of packets can be triggered by certain events, e.g., queue length reaching a pre-defined threshold, a timer having a time out, or some combination of these two.

Batching is typically accompanied by reordering. In this paper, the attacks focus on the traffic characteristics. As reordering does not change packet interarrival times much for mixes using batching, these attacks (and our analysis) are unaffected by reordering. Thus, our results are applicable to systems that use any kind of reordering methods. As such, in the rest of this paper, we will not discuss reordering techniques further.

Any of the batching strategies can be implemented in two ways:

Link-Based Batching: With this method, each output link has a separate queue. A newly arrived packet is put into a queue depending on its destination (and hence the link associated with the queue). Once a batch is ready from a particular queue (per the batching strategy), the packets are taken out of the queue and transmitted over the corresponding link.

Mix-Based Batching: In this way, the entire mix has only one queue. The selected batching strategy is applied to this queue. That is, once a batch is ready (per the batching strategy), the packets are taken out the queue and transmitted over links based on the packets' destination.

Each of these two methods has its own advantages and disadvantages. The control of link-based batching is distributed inside the mix and hence it may have good efficiency. On the other hand, mix-based batching uses only one queue and hence is easier to manage. We consider both methods in this paper.

3.3 Threat Model

In this paper, we assume that the adversary uses a classical timing analysis attack ([1, 22]), which we summarize as follows:

1. The adversary observes input and output links of a mix, collects the packet interarrival times, and analyzes them. This type of attack is passive, since traffic is not actively altered (by, say, dropping, inserting, and/or modifying packets during a communication session), and is therefore often difficult to detect. This type of attack can be easily staged on wired and wireless links [23] by a variety of agents, such as malicious ISPs or governments ([24, 25]).
2. To maximize the power of the adversary, we assume that she makes observations on all the links of the mix network.
3. The mix's infrastructure and strategies are known to the adversary. This is a typical assumption in the study of security systems. The above two assumptions create the worst case in terms of security analysis.
4. The adversary cannot correlate (based on packet timing, content, or size) a packet on an input link to another packet on the output link. Packet correlation based on packet timing is prevented by batching, and correlation based on content and packet size is prevented by encryption and packet padding, respectively.
5. To simplify the following discussion, we assume that dummy traffic is not used in the mix network. Some of the modern anonymous communication systems such as

Onion routing ([26]) do not use dummy traffic because of its heavy consumption of bandwidth and the general lack of understanding of to what extent exactly dummy packets contribute to anonymity.

6. Finally, we assume that the specific objective of the adversary is to identify the output link of a traffic flow that appears on an input link. Others have described similar attacks, but under simplified circumstances. Serjantov and Sewell [20], for example, assume that the flow under attack is alone on a link thus making its traffic characteristics immediately visible to the attacker. In this paper, we consider flows inside (potentially large) aggregates, thus making the attack generally applicable.

4 Traffic Flow Correlation Techniques

This section discusses the traffic flow correlation techniques that may be used by the adversary either to correlate senders and receivers directly or to greatly reduce the searching time for such a correlation in a mix network.

4.1 Overview

Recall that the adversary's objective is to correlate an incoming flow to an output link at a mix. We call this *flow correlation*. This kind of flow correlation attack is harmful in many scenarios. For example, in Figure 1, the adversary can discover the communication relationship between senders (S_1 and S_2) and receivers (R_1 and R_2) by matching senders' output flows and receivers' input flows. Using the flow correlation attack techniques, the adversary can find out a flow's sender and receiver if she catches a fragment of the flow in the mix network, thus breaking the anonymity despite the mix network. In a peer-to-peer mix network, the adversary can even reconstruct the path of this TCP connection by using flow correlation techniques. This subsection discusses the attack in more detail.

Figure 2 shows a flowchart of the typical procedure which the adversary may use to perform flow correlation. We now describe each step in detail.

(1) Data Collection. We assume that the adversary is able to collect information about all the packets on both input and output links. For each collected packet, the arrival time is recorded using tools such as tcpdump [27] and Cisco's NetFlow [28]. We assume that all the packets are encrypted and padded to the same size, and hence only arrival time is of interest. The arrival times of packets at input link i form a time series

$$A_i = (a_{i,1}, \dots, a_{i,n}) \quad (1)$$

where $a_{i,k}$ is the k^{th} packet's arrival time at input link i , and n is the size of the sample collected during a given sampling interval. Similarly, the arrival times of packets at output link j form a time series

$$B_j = (b_{j,1}, \dots, b_{j,m}) \quad (2)$$

where $b_{j,k}$ is the k^{th} packet's arrival time at output link j , and m is the size of the sample collected during a given sampling interval. The packets come out from mixes in batches.

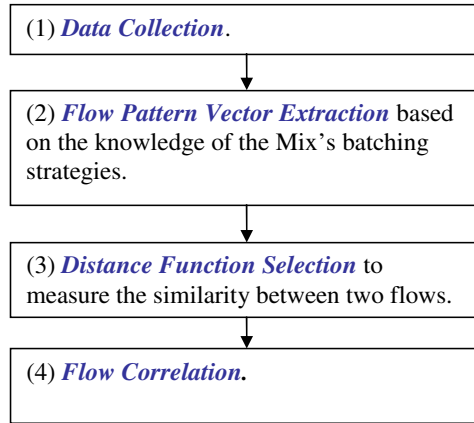


Fig. 2. Typical Flowchart for Flow Correlation

The length of sampling interval usually is much longer than the duration of a batch. Hence, a sampling interval typically contains many batches. We make the simplifying assumption that the traffic characteristic of the flow under consideration (the *input flow*) is known. This can be the case for example because the flow traffic characteristic is indeed observable at the input or at the input of the mix network.

(2) Flow Pattern Vector Extraction. With the above notation, the strategy of the adversary is to analyze the time series A_i s and B_j s in order to determine if there is any “similarity” between an input flow and an output flow of the mix. However, a direct analysis over these time series will not be effective. They need to be transformed into so called *pattern vectors* that can facilitate further analysis. We have found that effective transformation depends on batching strategies utilized by the mix. In Section 4.3, we will discuss specific definitions of transformations for different batching strategies. Currently, for the convenience of discussion, let us assume that A_i is transformed into pattern vector $X_i = (x_{i,1}, \dots, x_{i,q})$. And time series B_j is transformed into $Y_j = (y_{j,1}, \dots, y_{j,q})$. Note, here the two pattern vectors have the same length.

(3) Distance Function Selection. We define the distance function $d(X_i, Y_j)$, which measures the “distance” between an input flow at input link i and the traffic at output link j . The smaller the distance, the more likely the flow on an input link is correlated to the corresponding flow on the output link. Clearly, the definition of the distance function is the key in the correlation analysis. Section 4.2 will discuss two effective distance functions: one is based on mutual information and the other is based on the frequency-spectrum-based matched filter.

(4) Flow Correlation. Once the distance function has been defined between an input flow and an output link, we can easily carry out the correlation analysis by selecting the output link whose traffic has the minimum distance to input flow pattern vector X_i .

This approach can be easily extended to cases when multiple flows are aggregated over an input link [21]. The key idea is that by properly calculating the distance, we can find a correlation between one input flow and a set of output flows.

4.2 Flow Pattern Vector Extraction

In this subsection, we discuss how to choose pattern vectors X_i s and Y_j s. We will start with pattern vectors for the output link traffic first. Recall that batching strategies in Table 1 can be classified into two classes: threshold triggered batching (S_1 , S_3 , and S_5)¹ and timer triggered batching (S_2 , S_4 , S_6 and S_7). We will see that different classes should have different transformation methods.

For threshold triggered batching strategies, packets come out from the mix in batches. Hence, the inter-arrival time of packets in a batch is determined by the transmission latency, which is independent of the input flow. Thus, the useful information to the adversary is the number of packets in a batch and the time elapses between two batches. Normalizing this relationship, we define the elements in pattern vector Y_j as follows:

$$Y_{j,k} = \frac{\text{Number of packets in batch } k \text{ in the sampling interval}}{(\text{Ending time of batch } k) - (\text{Ending time of batch } k-1)} \quad (3)$$

In the calculation, we may need to truncate the original time series $B_j = (b_{j,1}, b_{j,2}, \dots, b_{j,n})$ so that only complete batches are used.

For timer triggered batching strategies, a batch of packets is sent whenever a timer fires. The length of the time interval between two consecutive timer events is a pre-defined constant. Thus, following a similar argument made for the threshold triggered batching strategies, we define the elements in pattern vector Y_j as follows:

$$Y_{j,k} = \frac{\text{Number of packets in the } k^{th} \text{ time out interval}}{(\text{time of } k^{th} \text{ time-out}) - (\text{time of } (k-1)^{st} \text{ time-out})} \quad (4)$$

$$= \frac{\text{Number of packets in the } k^{th} \text{ time out interval}}{\text{Pre-defined inter-time-out length}} \quad (5)$$

Again, in the calculation, we may need to truncate the original time series B_j so that only complete batches are used.

For the traffic *without batching* (i.e., the baseline strategy S_0 defined in Table 1), we use similar methods defined for timer triggered batching strategies as shown in (5).

The basic idea in the methods for extraction of pattern vectors is to partition a sampling interval into multiple sub-intervals and calculate the average traffic rate in each sub-interval as the values of the elements of traffic pattern vectors. The above two methods differ on how to partition the interval, depending on which batching strategy is used by the mix. We take a similar approach to extract pattern vectors X_i s corresponding to Y_j s. Again, the specific method of sub-interval partition depends on how the mix is batching the packets. Due to the space limitation, we will not further discuss the details of the methods developed. Readers are referred to [21] for details.

¹ S_3 could also be classified as timer-triggered. However, we treat it as threshold triggered because it may send out a batch when the number of packets received by the mix has reached the threshold.

4.3 Distance Functions

In this paper, we consider two kinds of distance functions: the first is based on a comparison of mutual information and the second on frequency analysis. The motivation and computation methods are given below.

Mutual Information

Mutual information is an information theoretical measure of the dependence of two random variables. In our scenario, we can view the pattern vectors that represent the input and output flows as samples of random variables. If we consider the pattern vectors X_i and Y_j to be each a sample of the random variables \mathcal{X}_i and \mathcal{Y}_j , respectively, then $\{(X_{i,1}, Y_{j,1}), \dots, (X_{i,q}, Y_{j,q})\}$ correspond to a sample of the joint random variable $(\mathcal{X}_i, \mathcal{Y}_j)$. With these definitions, the distance function $d(X_i, Y_j)$ between pattern vectors X_i and Y_j should be approximately inversely proportional to the mutual information $I(\mathcal{X}_i, \mathcal{Y}_j)$ between \mathcal{X}_i and \mathcal{Y}_j ,

$$d(X_i, Y_j) = \frac{1}{I(\mathcal{X}_i, \mathcal{Y}_j)} = - \frac{1}{\int \int p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(x_i)p(y_j)}} \quad (6)$$

Here, we need to estimate marginal distributions ($p(x_i)$ and $p(y_j)$) and their joint distribution $p(x_i, y_j)$. In this paper, we use histogram-based estimation of mutual information $\hat{I}(\mathcal{X}_i, \mathcal{Y}_j)$ of continuous distributions [29], which is given as follows.

$$\hat{I}(\mathcal{X}_i, \mathcal{Y}_j) \approx \sum_{u,v} \frac{K_{uv}}{q} \log \frac{K_{uv}N}{K_{u.}K_{.v}} \quad (7)$$

where q is the sample size. The sample space is a two-dimensional plane divided into $U \times V$ equally-sized $\Delta X \times \Delta Y$ cells with coordinates (u, v) . K_{uv} is the number of samples in the cell (u, v) . ΔX and ΔY have to be carefully chosen for an optimal estimation.

Frequency Analysis

For timer-triggered batching strategies, we therefore use FFT or Wavelet on the sample X_i and Y_j to obtain the frequency spectrum X_i^F and Y_j^F . Then we apply matched filter method over X_i^F and Y_j^F . We take advantage of the fact that frequency components of the input flow traffic carry on to the aggregate flow at the output link. Matched filter is an optimal filter to detect a signal buried in noise. It is optimal in the sense that it can provide the maximum signal-to-noise ratio at its output for a given signal. In particular, by directly applying the theory of matched filters, we can define the distance function $d(X_i, Y_j)$ as the inverse matched filter detector $M(X_i^F, Y_j^F)$,

$$d(X_i, Y_j) = \frac{1}{M(X_i^F, Y_j^F)} = \frac{1}{\frac{\langle X_i^F, Y_j^F \rangle}{\|Y_j^F\|}} \quad (8)$$

where $\langle X_i^F, Y_j^F \rangle$ is the inner product of X_i^F and Y_j^F , and $\|Y_j^F\| = \sqrt{\langle Y_j^F, Y_j^F \rangle}$. Please refer to [30] for details about the calculation of FFT and Wavelet over a vector. Due to the space limit, please refer to [21] for detailed results of the Wavelet-based method, which has similar results to the FFT method reported in this paper.

5 Empirical Evaluation

In this section, we evaluate the effectiveness of a selection of batching strategies (listed in Table 1) for a mix under our flow correlation attacks. We will see the failure of a mix under our traffic flow correlation attacks and batching strategies' influence on TCP flow performance.

5.1 Experiment Network Setup

Figure 3 shows our experimental network setup. Our mix is implemented on Timesys/Real Time Linux operating system for its timer accuracy [31]. The Mix control module that performs the batching and reordering functions is integrated into Linux's firewall system [32] using *Netfilter*; we use the corresponding firewall rules to specify what traffic should be protected. Two delay boxes D_1 and D_2 emulate the Internet propagation delay on different paths.

Our experiments reported here focus on TCP flows because of their dominance in the Internet. However, the results are generally applicable to other kinds of flows. The traffic flows in our experiments are configured as follows: An FTP client on node R_2 downloads a file from the FTP server on S_2 . The traffic from S_1 to R_2 serves as the random noise traffic to the FTP client. The traffic from node S_1 to node R_1 is the cross traffic through mix M from the perspective of the FTP flow. We maintain the traffic rate on both output links of the mix at approximately 500 packets per second (*pps*). The objective of the adversary in this experiment is to identify the output link that carries the FTP flow.

5.2 Metrics

We use *detection rate* as a measure of the ability of the mix to protect anonymity. Detection rate here is defined as the ratio of the number of correct detections to the number of attempts. While the detection rate measures the *effectiveness* of the mix, we measure its *efficiency* in terms of quality of service (QoS) perceived by the applications. We use *FTP goodput* as an indication of FTP quality of service (*QoS*). FTP goodput is defined as the rate at which the FTP client R_2 receives data from the FTP server S_2 . Low levels of FTP goodput indicate that the mix in the given configuration is poorly applicable for low-latency flow-based mix networks.

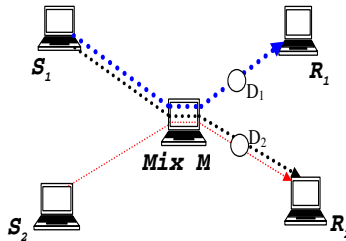
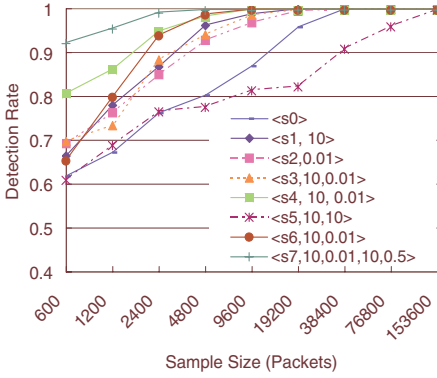
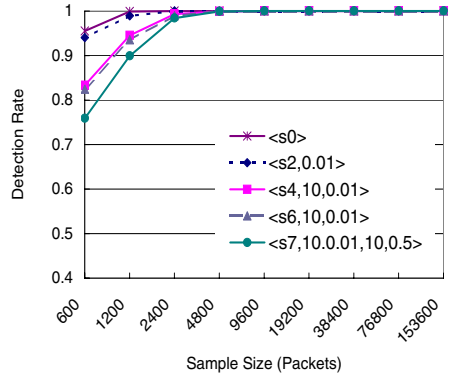


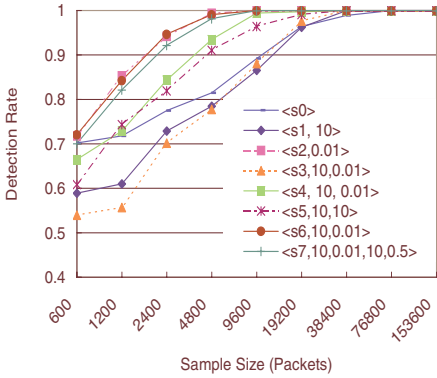
Fig. 3. Experiment Setup



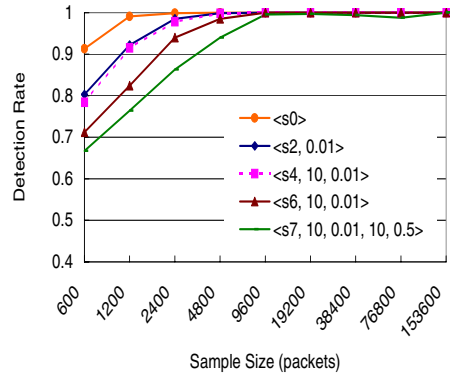
(a) Mutual Information



(b) Matched Filter

Fig. 4. Detection Rate for Link-based Batching

(a) Mutual Information



(b) Matched Filter

Fig. 5. Detection Rate for Mix-based Batching

5.3 Performance Evaluation

Effectiveness of Batching Strategies

Figure 4 shows the detection rate for systems using a link-based batching strategy. Figure 5 shows the detection rate for systems using a mix-based batching strategy as a function of the number of packets observed. A sample may include both FTP packets and cross traffic packets while FTP packets account for less than 20% of the number -sample size- of packets. Parameters in the legends of these figures are listed in the same order as in Table 1. Based on these results, we make the following observations:

1. For all the strategies, the detection rate monotonically increases with increasing amount of available data. The detection rate approaches 100% when the sample size is sufficiently large. This is consistent with intuition, as more data implies that there is more information about the input flow, which in turn improves the detection rate.

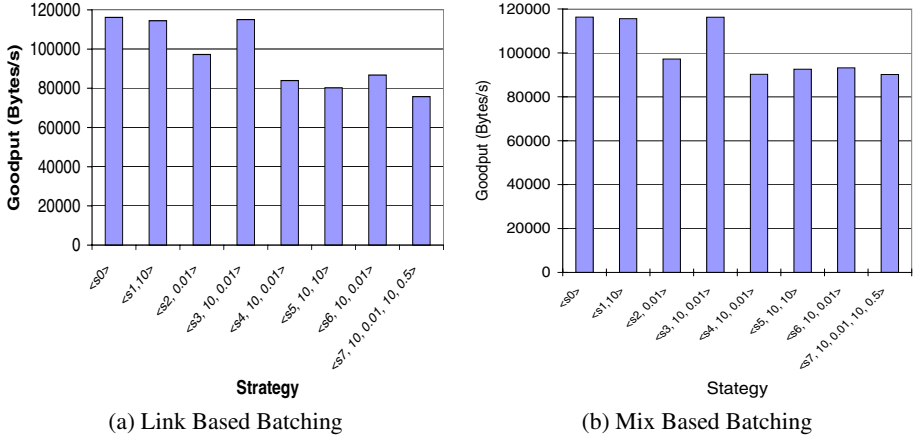


Fig. 6. FTP Goodput

2. Different strategies display different resistances to flow correlation attacks. In general, pool mixes perform better than simple mixes based on matched filter detector.
3. Frequency-analysis-based distance functions typically outperforms mutual-information-based distance functions in terms of detection rate. For many batching strategies, the former performs significantly better. This is because there are phasing issues in frequency-analysis-based attacks. Therefore, lack of synchronization between data collected at input and output port has a minor effect on the effectiveness of the attack.
4. To compare mix-based batching strategy with link-based batching strategy, we find that no one dominates the other.

Overall, our data shows that the mix using any of batching strategies S_1, S_2, \dots, S_7 fails under the flow correlation attacks. One of the reasons is that TCP flows often demonstrate interesting patterns such as periodicity of rate change and burstiness in particular when the TCP loop-control mechanism is triggered by excessive traffic perturbation in the mixes. Figure 4 and 5 show that flow correlation attacks can well explore this pattern difference between TCP flows.

Efficiency of Batching Strategies

As batching delays packets, one should expect that the overall performance (in terms of throughput) of TCP connections will be impacted by the mixes along their path. Figure 6 quantitatively shows the degradation of FTP goodput for a mix using different batching strategies.

In Figure 6, we compare FTP goodput between a strategy without any batching (S_0) and other batching strategies (S_1, S_2, \dots, S_7). We still use the network setup in Figure 3. The traffic other than FTP is configured as follows: 400pps from S_1 to R_1 and 500pps from S_2 to R_2 . Based on these experiments and the results illustrated in Figure 6, we make the following observations:

1. FTP goodput is decreased because of the use of batching.
2. Different batching strategies have different impact on the FTP goodput. In general, pool batching strategies (strategy S_5 to S_7) cause a worse FTP goodput than simple batching strategies (strategy S_1 to S_4).
3. When the batching in the mixes is excessively aggressive, that is, when batching intervals are too long or threshold values too high, the batching interferes with the time-out behavior of TCP and FTP, and in some cases, FTP aborts. This is the case in particular for threshold triggered mixes with no cross traffic.

6 A Countermeasure and Its Performance

From the discussion above, it is apparent that traditional batching strategies and reordering are not sufficient for mixes to effectively counter flow correlation attacks. Additional measures are needed. In this section, we introduce a relatively efficient and effective countermeasure and evaluate its performance in terms of FTP goodput.

6.1 Overview

A class of possible countermeasures can be developed based on the lessons learned in the previous sections. If a flow correlation attack relies on comparisons of pattern vectors of outgoing traffic, it will be ineffective when all packet vectors are identical. Thus, this type of flow correlation attacks can be effectively countered if a mix can make all the output flows look identical. As a result, assuming that we have the input flow vector X_i and l output flow vectors Y_1, \dots, Y_l ,

$$d(X_i, Y_1) = \dots = d(X_i, Y_j) = \dots = d(X_i, Y_l), \quad (9)$$

and the only analysis strategy for an adversary would be to randomly guess which output flow is correlated to an input flow. This results in a detection rate of $\frac{1}{l}$.

Because naturally the rates of traffic along all the output links of a mix are different, we have to appropriately insert dummy packets to make all the output flows behave in the same way. A challenge here is to insert a minimum number of dummy packets.

Such an output-control algorithm is illustrated in Figure 7. Mix M maintains two output queues, Q_1 for the link between Mix M and node R_1 , and Q_2 for the link between

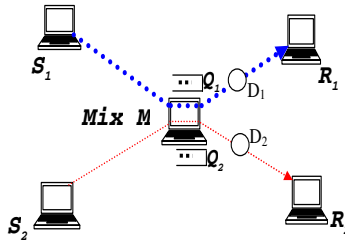


Fig. 7. Network Setup for the New Countermeasure

```

Data      : queues, in which packets are kept in deadline order by the mix
Result    : synchronized flows out of the mix
while (1) do
  if ( $Q_1.Length > 0$ ) and ( $Q_2.Length > 0$ ) then
    send the first packet from  $Q_1$ ;
    send the first packet from  $Q_2$ 
  else
    if ( $Q_1.Length > 0$ ) then
      if ( $Q_1.FirstPacket.Deadline > CurrentTime$ ) or ( $Q_1.Length > Q_1.Threshold$ )
        then
          send the first packet from  $Q_1$ ;
          send a dummy packet for  $Q_2$ 
        end
      end
    else
      if ( $Q_2.Length > 0$ ) then
        if ( $Q_2.FirstPacket.Deadline > CurrentTime$ ) or ( $Q_2.Length > Q_2.Threshold$ ) then
          send a dummy packet for  $Q_1$ 
          send the first packet from  $Q_2$ ;
        end
      end
    end
  end
end

```

Fig. 8. Algorithm for Output Traffic Control

Mix M and node R_2 . At any time, if each queue has a packet, they are sent out in some pre-defined order, e.g., the packet in Q_1 first and the packet in Q_2 second. By doing so, one of the two queues will be always empty. Let us say, for the moment, that Q_2 is empty. A deadline is assigned to each packet waiting in Q_1 . If a packet in Q_1 reaches its deadline, a dummy packet will be generated for Q_2 . Then, the payload packet from Q_1 and the dummy packet from Q_2 are sent out in the predefined order. A dummy packet will also be generated for Q_2 if the queue length of Q_1 goes beyond a preset threshold. In this way, we can ensure a maximum delay on each packet, and we also guarantee that neither queue will overflow.

Figure 8 gives the new countermeasure algorithm on Mix M for the anonymity system in Figure 7. We can see that the output traffic of the Mix is now synchronized, and the adversary cannot observe any difference among the output flows.

This method can be easily extended and optimized for more complicated cases. The number of virtual output links of a mix can be very large since we assume a peer-to-peer mix network. Since we only maintain virtual queues, the overhead is limited. In the case of a large network with a small number of flows, there still needs to be a lower bound LB_Q of the number of virtual queues required for each mix to maintain anonymity. In other words, we do not necessarily need to synchronize every output link when traffic is slow, but we will synchronize a minimum number LB_Q of links. For example, if there is one virtual queue with a packet whose deadline is reached, we have to send out dummy packets to the other $LB_Q - 1$ virtual links.

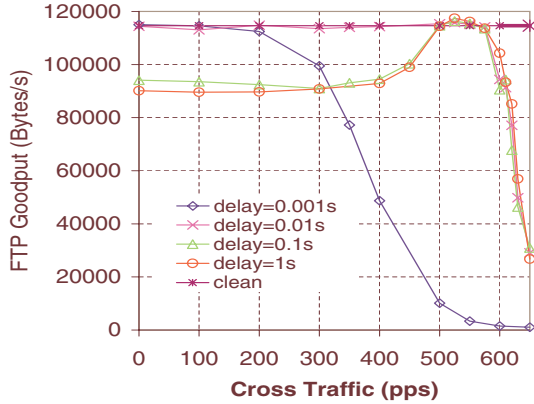


Fig. 9. FTP Goodput Using Output Traffic Control (“clean” means no output traffic control)

Output traffic control is not new and has been proposed for example in [33], where messages at the output ports are forwarded periodically². The algorithm in Figure 8 is more efficient and probably more effective than the approach described in [33]. It is more efficient because packets are forwarded based on each queue’s status: once each queue has payload packets, the first packet in each queue is sent out and packets suffer smaller delay at Mixes. It is likely more effective because periodic traffic patterns are very difficult to generate with sufficient accuracy. We showed in NetCamo [22, 34], for example, how high-accuracy traffic analysis can easily break periodic link padding schemes.

6.2 Performance Evaluation of Output Traffic Control

We are interested in how traffic flows traversing a mix affect each other. In particular, we evaluate the TCP performance. Again FTP is used as an example in the evaluation.

Figure 9 gives the FTP goodput measurement for our new scheme for the network setup in Figure 7. We set the threshold of each queue at 50 packets. The path from S_2 to R_2 has FTP traffic and UDP traffic of 400pps. Cross traffic in Figure 9 refers to the UDP traffic along the path S_1 to R_1 . Both paths have a propagation delay of 0.3 second. We have the following observations from these experiments:

1. While not evident from Figure 9, the observed detection rate of the correlation attack is 50% in all the cases when the new countermeasure is used. This is expected, as the new method can guarantee a detection rate of $1/LB_Q$ where $LB_Q = 2$ in this case.
2. The goodput for the clean FTP is 114,628.83 bytes/s. When the delay parameter is set to 0.01s, the same goodput is achieved as long as the cross traffic is less than 525 pps. This is very significant. It indicates that, once the delay parameter is properly

² The paper is too vaguely written for us to figure out exactly what forwarding mechanism is used.

selected, our new method can achieve high throughput (as high as the case without mix) while guaranteeing a low detection rate.

3. For the cases of delay equal to 0.01s, 0.10s, and 1.00s, right after the cross traffic goes beyond 525 pps, all have their goodput drop rapidly. This is due to the fact that the cross traffic is so heavy that the FTP's TCP protocol detects congestion and adapts accordingly.
4. It is also interesting to note, that when the cross traffic is low and the value of delay parameter is large (say, the cross traffic is less than 500 pps and delay is equal to 0.10s or 1.00s), the goodput is low (about 93,000 bytes/s). This is consistent with intuition: if the cross traffic is low and delay is large, then the traffic of our FTP flow may have to wait longer than in other cases, resulting in a reduction of goodput.
5. Finally, in the case when the value of delay parameter is small, say, equal to 0.001s, the curve of goodput is monotonically decreasing. In this case, it is likely that a packet from the FTP flow will be transmitted due to the deadline expiration, rather than the arrival of a packet from the cross traffic. Thus, the cross traffic always contributes negatively to the goodput performance here by creating dummy packets.

7 Summary and Future Work

We have analyzed mix networks in terms of their effectiveness in providing anonymity and quality-of-service. Various methods used in mix networks were considered: seven different packet batching strategies and two implementation schemes, namely the link-based batching scheme and mix-based batching scheme. We found that mix networks that use traditional batching strategies, regardless of the implementation scheme, are vulnerable under flow correlation attacks. By using proper statistical analysis, an adversary can accurately determine the output link used by traffic that comes to an input flow of a mix. The detection rate can be as high as 100% as long as enough data is available. This is true even if heavy cross traffic exists. The data collected in this paper should give designers guidelines for the development and operation of mix networks.

The failure of traditional mix batching strategies directly leads us to the formation of a new packet control method for mixes in order to overcome their vulnerability to flow correlation attacks. Our new method can achieve a guaranteed low detection rate while maintaining high throughput for normal payload traffic. Our claim is validated by extensive performance data collected from experiments. The new method is flexible in controlling the overhead by adjusting the maximum packet delay.

Our study is the first that systematically models and analyzes flow correlation attacks and their countermeasures. The work presented in this paper is largely empirical. We are currently developing an analysis framework that allows quick, back-of-the-envelope calculations to assess the effectiveness of batching strategies in countering flow correlation attacks. It is an open question what statistical analysis methods an adversary may use. Performance bounds and estimates in terms of detection rate and throughput may be developed by following the approaches taken in [35] and [36], respectively.

References

1. Song, D.X., Wagner, D., Tian, X.: Timing analysis of keystrokes and timing attacks on ssh. In: Proceedings of 10th USENIX Security Symposium. (2001)
2. Sun, Q., Simon, D.R., Wang, Y., Russell, W., Padmanabhan, V.N., Qiu, L.: Statistical identification of encrypted web browsing traffic. In: Proceedings of IEEE Symposium on Security and Privacy. (2002)
3. Chaum, D.: Untraceable electronic mail, return addresses, and digital pseudonyms. *Communications of the ACM* **4** (1981)
4. Helsingius, J.: Press release: Johan helsingius closes his internet remailer. <http://www.penet.fi/press-english.html> (1996)
5. Parekh, S.: Prospects for remailers - where is anonymity heading on the internet. <http://www.firstmonday.dk/issues/issue2/remailers/> (1996)
6. Gülcü, C., Tsudik, G.: Mixing E-mail with Babel. In: Proceedings of the Network and Distributed Security Symposium (NDSS), IEEE (1996) 2–16
7. Möller, U., Cottrell, L.: Mixmaster Protocol — Version 2. <http://www.eskimo.com/~rowdenw/crypt/Mix/draft-moeller-mixmaster2-proto\%col-00.txt> (2000)
8. Serjantov, A., Dingledine, R., Syverson, P.: From a trickle to a flood: active attacks on several mix types. In: Proceedings of Information Hiding Workshop. (2002)
9. Danezis, G., Dingledine, R., Mathewson, N.: Mixminion: Design of a Type III Anonymous Remailer Protocol. In: Proceedings of the 2003 IEEE Symposium on Security and Privacy. (2003)
10. Back, A., Möller, U., Stiglic, A.: Traffic analysis attacks and trade-offs in anonymity providing systems. In: Proceedings of Information Hiding Workshop. (2001)
11. Berthold, O., Langos, H.: Dummy traffic against long term intersection attacks. Proceedings of Privacy Enhancing Technologies workshop (PET 2002) (2002)
12. Berthold, O., Pfitzmann, A., Standtke, R.: The disadvantages of free MIX routes and how to overcome them. In: Proceedings of Designing Privacy Enhancing Technologies: Workshop on Design Issues in Anonymity and Unobservability. (2000)
13. Mitomo, M., Kurosawa, K.: Attack for Flash MIX. In: Proceedings of ASIACRYPT. (2000)
14. Raymond, J.: Traffic analysis: Protocols, attacks, design issues and open problems. In: Proceedings of Designing Privacy Enhancing Technologies: Workshop on Design Issues in Anonymity and Unobservability. (2001)
15. Syverson, P.F., Goldschlag, D.M., Reed, M.G.: Anonymous connections and onion routing. In: Proceedings of IEEE Symposium on Security and Privacy, Oakland, California (1997) 44–54
16. Boucher, P., Shostack, A., Goldberg, I.: Freedom systems 2.0 architecture (2000)
17. Reiter, M., Rubin, A.: Crowds: Anonymity for web transactions. *ACM Transactions on Information and System Security* **1** (1998)
18. Freedman, M.J., Morris, R.: Tarzan: A peer-to-peer anonymizing network layer. In: Proceedings of the 9th ACM Conference on Computer and Communications Security (CCS), Washington, DC (2002)
19. Sherwood, R., Bhattacharjee, B., Srinivasan, A.: p^5 : A protocol for scalable anonymous communication. In: Proceedings of the 2002 IEEE Symposium on Security and Privacy. (2002)
20. Serjantov, A., Sewell, P.: Passive attack analysis for connection-based anonymity systems. In: Proceedings of European Symposium on Research in Computer Security (ESORICS). (2003)
21. Zhu, Y., Fu, X., Graham, B., Bettati, R., Zhao, W.: Correlation attacks in a mix network. Texas A&M University Computer Science Technical Report TR2003-8-9 (2003)

22. Fu, X., Graham, B., Xuan, D., Bettati, R., Zhao, W.: Analytical and empirical analysis of countermeasures to traffic analysis attacks. In: *Proceedings of International Conference on Parallel Processing (ICPP)*. (2003)
23. Howard, J.D.: An analysis of security incidents on the internet 1989 - 1995. Technical report, Carnegie Mellon University Dissertation (1997)
24. F.B.I.: Carnivore diagnostic tool. <http://www.fbi.gov/hq/lab/carnivore/carnivore2.htm> (2003)
25. Walton, G.: China's golden shield: Corporations and the development of surveillance technology in china. <http://www.totse.com/en/privacy/privacy/pucc.html> (2003)
26. Achives, O.R.D.: Link padding and the intersection attack. <http://archives.seul.org/or/dev> (2002)
27. tcpdump.org: tcpdump. <http://www.tcpdump.org/> (2003)
28. cisco.inc.: Netflow services solutions guide. <http://www.cisco.com/univercd/cc/td/doc/cisintwk/intsolns/netflsol/nfwhite.htm> (2003)
29. Moddemeijer, R.: On estimation of entropy and mutual information of continuous distributions. *Signal Processing* **16** (1989) 233–246
30. MathWorks: Documentation for mathworks products (release 13). <http://www.mathworks.com/access/helpdesk/help/helpdesk.shtml> (2003)
31. TimeSys: Timesys linux docs. <http://www.timesys.com/> (2003)
32. netfilter.org: Netfilter. <http://netfilter.samba.org/> (2003)
33. Rennhard, M., Rafaeli, S., Mathy, L., Plattner, B., Hutchison, D.: Analysis of an anonymity network for web browsing. In: *Proceedings of Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises (WET ICE)*. (2002)
34. Guan, Y., Fu, X., Xuan, D., Shenoy, P.U., Bettati, R., Zhao, W.: Netcamo: Camouflaging network traffic for qos-guaranteed critical allplications. *IEEE Transactions on Systems, Man, and Cybernetics Part A: Systems and Humans, Special Issue on Information Assurance* **31** (2001) 253–265
35. Fu, X., Graham, B., Bettati, R., Zhao, W.: On effectiveness of link padding for statistical traffic analysis attacks. In: *Proceedings of the 23rd IEEE International Conference on Distributed Computing Systems (ICDCS 2003)*. (2003)
36. Padhye, J., Firoiu, V., Towsley, D., Krusoe, J.: Modeling TCP throughput: A simple model and its empirical validation. In: *Proceedings of ACM SIGCOMM Special Interest Group on Data Communications (SIGCOMM)*. (1998)