

SDA-2H: Understanding the Value of Background Cover Against Statistical Disclosure

Mahdi N. Al-Ameen, Charles Gatz, Matthew Wright

Dept. of Computer Science and Engineering, University of Texas at Arlington, USA
 mahdi.al-ameen@mavs.uta.edu, charles.gatz@mavs.uta.edu, mwright@cse.uta.edu

ABSTRACT

The statistical disclosure attack (SDA) is an effective method for compromising the anonymity of users in a mix based system. Cover traffic, in the form of fake or dummy messages sent by other users of the mix, is an effective defense to make the task of the attacker difficult. Our aim is to examine the effect that background cover – the cover traffic sent by other users – has on the effectiveness of statistical disclosure attacks. Since the original SDA does not explicitly account for background traffic volumes, we developed an extension to the SDA called SDA-2H that uses this information to improve upon the SDA. Based on this attack, we are able to quantify the importance of background cover traffic, which we show in simulation to be effective in many scenarios.

Keywords

online privacy, anonymous communications, statistical disclosure attack, background cover.

1. INTRODUCTION

To further ensure anonymity in online communication, Chaum introduced the concept of mixes in 1981 [1]. A mix can be understood as a node that acts as an intermediary between the sender and the receiver. It processes each packet of data before sending it to the recipient. This processing includes encrypting the data with the mix's public key (to be decrypted when it leaves the mix) as well as reordering the timing of the data packets (e.g. delaying the data packets) to obfuscate the path from sender to receiver. Using a set of mixes, a user can effectively stop an external eavesdropper from linking a receiver to a particular sender [10].

Despite the obvious advantages of mix networks in preserving anonymity, they are still vulnerable to a set of attacks. One set of such attacks are the intersection attacks (also referred to as partitioning or disclosure attacks). These perform traffic analysis by exploiting the fact that different messages use the same route through the network. They are particularly devastating because they do not rely upon any particular properties of mixing [2].

To mount these attacks, a passive eavesdropper observes a large volume of network traffic and relies on the assumption that a targeted sender (henceforth referred to as Alice) sends messages only to a limited number of recipients. Intersection attacks are designed based on correlating the times when senders and receivers are active. By observing which recipients received packets during the rounds when Alice was sending, the attacker can create a list of Alice's most frequent recipients, thereby thwarting her attempt at anonymity.

Though effective, the disclosure attack is computationally difficult and time consuming. The statistical disclosure attack (SDA) can be considered as an improvement on disclosure attacks. Unlike disclosure attacks, which involve heavy computation, SDA simply relies on performing simple operations over many observations [2,9]. The main question is how many rounds of observations are needed for the attacker to learn about Alice's sending behavior.

Users may attempt to make traffic analysis attacks like the SDA harder to perform by sending cover traffic, which consists of fake (dummy) messages that appear to be real messages and mask Alice's true sending activity. Mathewson and Dingledine showed the cover traffic is an effective defense against the SDA [9]. Perhaps surprisingly, however, the cover traffic that other users send has minimal effect on Alice's anonymity [7]. One might expect to find a difference between a scenario in which other senders greatly vary their cover traffic volume and when other senders do not use any cover traffic at all. The reason that it does not have an effect is that the SDA equation does not consider background traffic volume at all.

In this paper, we examine the value of background cover traffic as a defense against SDA. In particular, we extend the SDA into a modified attack that we call SDA-2H (SDA with Two Heads). In the SDA-2H, background traffic volumes are used. We have discovered that background traffic volumes can be used to estimate the amount of cover traffic that Alice sends, separately from the number of real messages that she sends. This allows the attacker to better estimate Alice's activity.

The paper is organized as follows: (§2) provides background on SDA, mix types, and cover traffic. (§3) explains our attacker model, SDA-2H. (§4) details the implementation and testing of our hypothesis, with the results explained and compared in (§5). Related work is mentioned in (§6), and our thoughts on future work are addressed in (§7). Finally, in (§8) we summarize our contributions.

2. MIX NETWORKS AND THE STATISTICAL DISCLOSURE ATTACK

In this section, we give a brief background on mix networks, before describing the SDA in detail and the different types of cover traffic that can be used to combat traffic analysis.

2.1 Mixes

Mix networks are a common choice for implementing anonymous systems and remain an active area of research. As explained in [1] and our introduction, a mix is a very powerful tool for obtaining anonymity — though it is bound by a few constraints [10]. There are various kinds of mix models available; however, in this project we are concerned with only binomial mixes.

In the binomial mix model, the decision to either send the mes-

sage in the current round, or delay until a later round is made by subjecting each incoming message to a coin toss with a bias probability, $pDelay$. As $pDelay$ approaches one, more messages are held in the mix, which makes correlating ingoing and outgoing messages difficult.

2.2 Statistical Disclosure Attack

SDA can be considered as an extension of intersection attacks, which function by analyzing the pattern of messages being sent from a particular user. These attacks are based on the fact that different messages are sent through the mix-based network using the same path. SDA takes this further by assuming that a user (Alice) sends messages only to a limited set of recipients. By recording and comparing sending and receiving patterns for each round, an attacker can estimate Alice’s most likely recipients by storing the cumulative probability that Alice sends to any given recipient. SDA is based on long term observations and consists of relatively simple mathematical processing, thus it is cheap and easy to implement. The accuracy of SDA depends on the number of rounds observed, weighed against any defensive strategies in place (*i.e.*, more observations yield a greater chance of noticing a difference between Alice’s sending pattern and the background; the use of cover traffic aims to make this period of time unacceptable to a prospective attacker by introducing error into the equation).

As mentioned in the introduction, SDA can also be considered as an extension to the disclosure attack described by Kesdogan [6]. The SDA equation is given as follows:

$$\bar{O} = \frac{\bar{m}\bar{v} + (\bar{n} - \bar{m})\bar{u}}{\bar{n}} \quad (1)$$

where \bar{O} is the mean of \vec{o} , a vector which records the probability that a given message was sent to each recipient; \bar{m} and \bar{n} are the average number of messages sent by Alice, and the average total number of messages sent (including Alice), in a round, respectively. The unknown sending behavior of Alice is denoted \bar{v} , and the likelihood that the background (everyone but Alice) sends to recipients, in the vector \bar{u} (where each element \bar{v}_i is the likelihood that Alice sends to recipient i and similarly for \bar{u}).

2.3 Cover Traffic

The cover traffic in a mix-based anonymity system consists of dummy messages that are added to the network along with the real messages transmitted by users of the network. Dummy messages serve as a useful tool to increase anonymity, since encryption inside the mix makes them indistinguishable from real messages (usually, until they are sent). The only option an attacker has is to account for all the messages into his observations, thus increasing the total number of messages sent by the user (which should make his work more difficult). Dummy messages can be classified into three types based on their origin:

- **User cover:** generated by the user Alice.
- **Background cover:** generated by senders other than Alice in the system.
- **Receiver-bound cover:** generated by the mix.

As our project centers on background cover, we now explain it in more detail.

2.4 Background Cover

Background cover is created when users (not part of the mix) generate dummies along with their real messages. As shown in [7],

this can be very effective in confusing a naive attacker. The attacker sees more messages entering the mix, and expecting them to exit the mix, makes a miscalculation. This demands more rounds of observation, and doesn’t threaten Alice’s privacy in the near future. The protection provided by background cover, however, can be easily nullified by a well-informed attacker, who can simply estimate the average number of dummy messages, \bar{d} (background cover) based on a simple equation:

$$\bar{d} \approx \bar{in} - \bar{out} \quad (2)$$

where \bar{in} is the mean number of messages entering the mix, and \bar{out} , the mean number of messages exiting the mix. Accounting for the percentage of background cover in their calculations; an attacker can virtually disregard the effect of background cover. It is proposed in [7] that an informed attacker can accurately estimate this percentage, thus reducing a receiver’s anonymity.

3. ATTACKER MODEL (SDA-2H)

We now describe SDA with Two Heads (SDA-2H), an extension of the SDA that seeks to extract more information from the same set of observations that the SDA uses. The key intuition of SDA-2H is that the attacker can estimate how much of Alice’s traffic is cover traffic based on the difference between the volume of incoming traffic and the volume of outgoing traffic. Prior work has shown that Alice cover provides an important protection against the [7,9]. If the attacker knows how much of Alice’s traffic is cover traffic, the SDA can operate perfectly, as if Alice was not sending cover traffic at all.

Let us illustrate this with a simple numerical example. Suppose that, in a given round, the background senders transmit 100 messages (all real messages) and Alice sends ten messages. Thus 110 messages reach the mix. The attacker observes that 105 messages leave the mix in that round. Since the attacker knows that the difference between the number of messages entering and leaving the mix is five, he can easily infer that Alice sent five cover messages and five real messages. Thus, the attacker can remove the effect of Alice cover.

More generally, the attacker needs to obtain \bar{O} to use in Equation 1 and does so by knowing the delay policy and observing the number of messages entering and exiting the mix. The number of Alice’s messages exiting the mix is denoted as n_{Alice} , with the other outgoing messages stored as $n_{Background}$. The set of recipients is represented by \vec{r} , which contains an element for every one of the mix’s recipients such that $\vec{r}[i]$ is a count of the messages received by the i^{th} recipient on a given round. This allows for updating the value of \bar{O} as follows:

$$\bar{O}[i] = \frac{\vec{r}[i] * n_{Alice}}{n_{Alice} + n_{Background}} \quad (3)$$

However, Malleh et al. showed that dummies sent into the mix can be calculated by an informed attacker [7]. The improvement that SDA-2H provides is to consider that if the actual background traffic is known (say B_{Real}), and Alice’s ingoing messages to the mix are represented as A_{Total} , then in , (for Equation 2), can be easily derived as $in = A_{Total} + B_{Real}$. At the same time, Alice’s cover is computed as $A_{Cover} = in - out$, which combines to derive Alice’s actual number of messages, $A_{Real} = A_{Total} - A_{Cover}$. The more accurate value of A_{Real} is then used in place of n_{Alice} in Equation 3.

4. SIMULATION

To evaluate the effectiveness of SDA-2H, we have extended the simulation environment used by Nayantara Malleš in her investigation of SDA. We refer the reader to her papers for additional details [7,8].

The basic model of the simulation is of a single mix that receives and sends messages in rounds. In each round, the user of interest (Alice) sends message to a subset of her contacts; the number of messages is selected from a geometric distribution. She distributes messages to her contacts evenly. The other users (the background) send messages according to a normal distribution with a standard deviation of 10% of the mean number of messages sent per round. The senders select recipients according to a scale-free model, chosen to be similar to a social network, with more popular receivers getting more messages.

In our simulation, we study binomial mixes. In a sense, a binomial mix can be seen as a rough approximation of a pool mix, following the work of both Malleš and Wright [7] and Mathewson and Dingledine [9] (both of whom describe their simulations as being on a pool mix, which was not correct [8]).

We have simulated for different values of background messages, that follow a normal distribution with mean in the range from 100 to 10000. We have simulated for dynamic background cover that vary in the ranges of [1–10]%, [45–55]%, [80–120]% and [160–240]%. The percentages of non-dynamic background cover that we have simulated for: [10%, 25%, 50%, 100%, 200%, 300%].

Average real messages per round from Alice (Geometric) are varied through choosing the values: [0.1, 0.3, 0.6, 0.9] and this value is 0.6 for the experiments when it remains fixed. We have chosen the values: [0.1, 0.3, 0.6, 0.9] for varying average dummy messages per round from Alice ($pDummy$) and 0.6 for the experiments where $pDummy$ remains constant.

Our primary metric is the number of rounds until the attacker is successful. We generally say that the attacker is successful, and stop a given run of the simulation, when he correctly identifies 25% of Alice’s recipients (when total Alice’s recipients is greater than 20). For Alice’s recipients in the range between 10 to 20, the attacker has to identify 50% of Alice’s recipients and for total Alice’s recipients less than 10, the attacker gets successful if it successfully identifies all the recipients of Alice. While this is an arbitrary amount, it provides a good indicator of when the attacker has learned a substantial amount of Alice’s sending behavior. In measuring the attacker’s success, we define a *granularity* of sampling. Rather than checking to see if the attacker is successful every round, we check every ten rounds. Sometimes, the attacker is not successful for a very long time. We set a cutoff value of $Max = 10^6$ rounds, after which we stop the simulation and say that the attacker took 10^6 rounds, which is effectively a failure for most realistic scenarios.

5. RESULTS

We now present the results of our simulations.

Our results show that when the value of background messages(mean) is greater than 200, the attacker gets success only for the ranges of [1–10]% and [45–55]% for dynamic background cover with attacker adjustment. The attacker fails for all the simulated ranges of dynamic background cover when the value of background messages(mean) is greater than or equal to 1000. But when the value of background messages(mean) is less than or equal to 200, the attacker gets success for all the simulated ranges of dynamic background cover with attacker adjustment. In these cases, the number of Alice’s recipients is 64 and the number of total re-

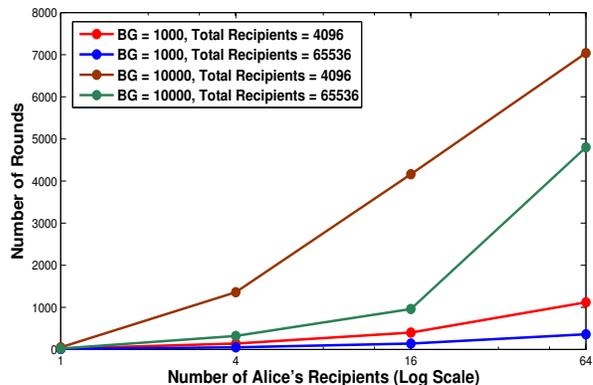


Figure 1: Number of Rounds [Non-Dynamic BG Cover(100%), Binomial mix]

cipients is 65536.

We have also simulated for different values of Alice’s recipients: [1, 4, 16, 64] and total recipients: [16, 4096, 65536]. We have found that for [Alice’s Recipients, Total Recipients] = [1, 4096] and [1, 65536], the attacker gets success for the ranges of [1–10]% and [45–55]%, when the background messages(mean) is 1000. The attacker is also successful for the same ranges of dynamic background cover, when the background messages(mean) is increased to 1700 for [Alice’s Recipients, Total Recipients] = [1, 65536]. But for background messages of 9000 and 10000 the attacker fails for any combination of Alice’s recipients and total recipients.

For dynamic background cover without attacker adjustment the attacker gets success only for the range of [1–10]%. When background messages(mean) is 100 the attacker gets success for all the simulated combinations of Alice’s recipients and total recipients for this range. When the value of background messages(mean) is increased to 1000, the attacker only gets success for [Alice’s recipients, Total recipients] = [1, 4096] and [1, 65536]. For background messages(mean), greater than or equal to 1700, the attacker fails to identify Alice’s recipients for dynamic background cover without attacker adjustment.

For non-dynamic background cover with attacker adjustment, the attacker successfully identifies Alice’s recipients for any simulated values of background messages. For background messages(mean) of 125, the attacker gets success for 10% of non-dynamic background cover without attacker adjustment. But when the value of background messages (mean) is greater than or equal to 1000, the attacker fails for any simulated value of non-dynamic background cover without attacker-adjustment.

5.1 Effects of Varying Parameters

Fig. 1 focuses on the effect of varying Alice’s recipients and total recipients with the change in background messages(mean) when the percentage of non-dynamic background cover is 100% (with attacker-adjustment). Here we find that for a fixed number of total recipients, the number of rounds, the attacker needs to be succeeded increases with the increase in the number of Alice’s recipients. The increases in the background messages makes it difficult for the attacker to be successful having total recipients and Alice’s recipients unchanged. Our results show that when Alice’s recipients remain unchanged and the number of total recipients increases, the attacker needs less number of rounds indeed, to successfully identify Alice’s recipients.

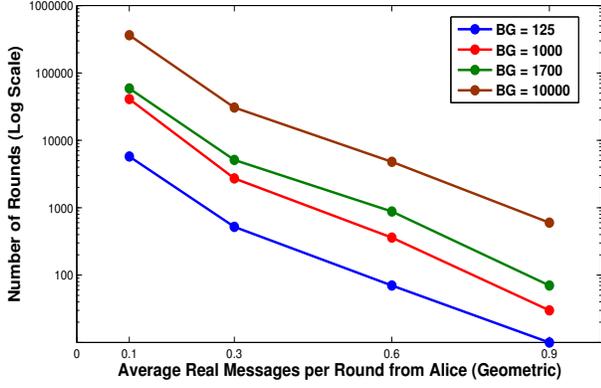


Figure 2: Number of Rounds [Non-Dynamic BG Cover(100%), Binomial mix]

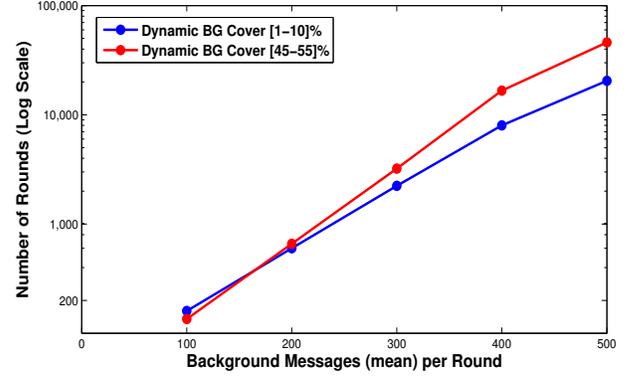


Figure 4: Number of Rounds [Binomial mix]

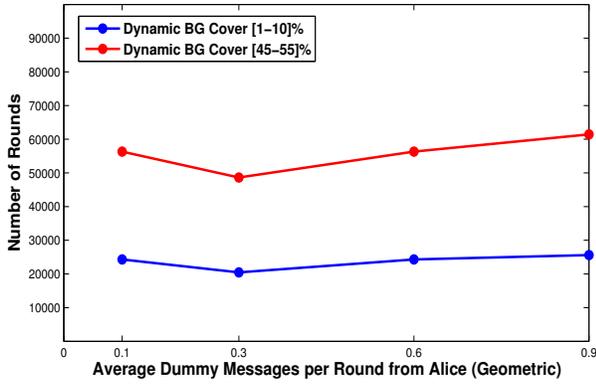


Figure 3: Number of Rounds [BG = 500, Binomial mix]

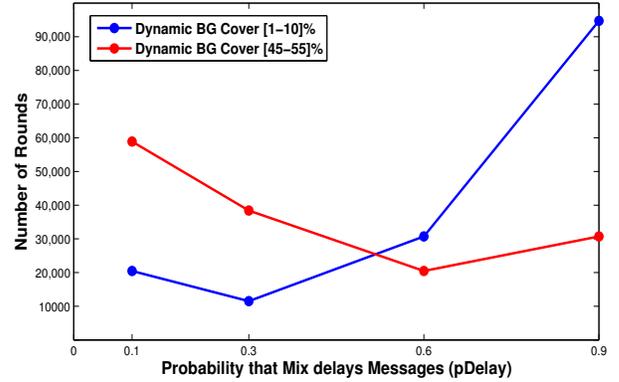


Figure 5: Number of Rounds [BG = 500, Binomial mix]

In Fig. 2, we get the results for varying average real messages per round from Alice with the change in background messages(mean), when the percentage of non-dynamic background cover is 100% (with attacker adjustment). From Fig. 2, we find that with the increase in the average real messages per round from Alice, the attacker needs less number of rounds to be successful. But with the increase in the background messages(mean), the number of rounds increases.

Fig. 3 reflects the effectiveness of SDA-2H attack model that reduces the effect of dummy messages per round from Alice (Alice Cover). From the Figure we find that for [1 – 10]% and [45 – 55]% of dynamic background cover with attacker adjustment, varying Alice Cover does not have significant impact on the number of rounds taken by the attacker to be successful.

From Fig. 4 we find the effect of varying background messages (mean) when the ranges of dynamic background cover (with attacker adjustment) are [1 – 10]% and [45 – 55]%. The results show that the number of rounds for the attacker to get success increases with the increase in background messages. The number of rounds for dynamic background cover with the range [45 – 55]% is less than [1 – 10]% when the value of background messages(mean) is 100 but with the increases in background messages the attacker needs more number of rounds for [45 – 55]% of dynamic background cover in comparison to that of [1 – 10]%. Fig. 5 focuses the effect of varying $pDelay$. Here we find that

for dynamic background cover varying in the range [1 – 10]%, the number of rounds increases with the increase in $pDelay$ from 0.3 to 0.9. The number of rounds is same for both ranges of [1 – 10]% and [45 – 55]% when the value of $pDelay$ is around 0.5. For dynamic background cover with the range of [45 – 55]%, the number of rounds decreases with the increase in $pDelay$ from 0.1 to 0.6 and an increase in number of rounds is found when $pDelay$ is increased from 0.6 to 0.9.

5.2 Comparison with Simple SDA

Simple SDA does not reduce the effect of Alice cover. So with the increase in Alice cover, number of rounds also increases significantly. But SDA-2H reduces the effect of Alice cover. So changing Alice cover does not have significant impact on the number of rounds in this case. In simulations without any background cover, our results stand for the fact that SDA-2H is more powerful than simple SDA and the results hold good for non-dynamic background cover with attacker adjustment(Fig. 6). But for dynamic background cover, the number of rounds taken by simple SDA is less than that of SDA-2H (Fig. 7). Our results show that in case of dynamic background cover, for $pDelay = 0.9$, SDA-2H is more powerful than simple SDA.

Once the amount and variability of background cover traffic are sufficiently large, it becomes a difficult task to correctly estimate the amount of Alice cover. At this point, SDA-2H with our simple averaging method for estimating Alice cover becomes error prone

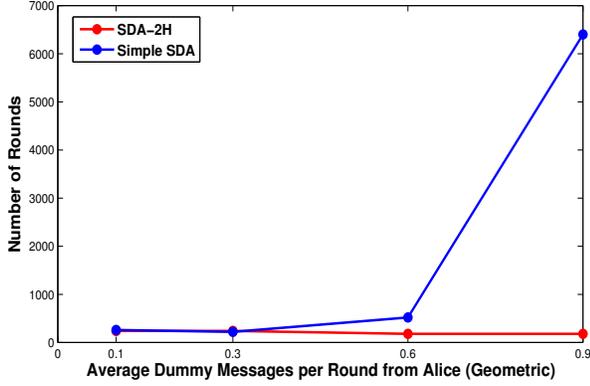


Figure 6: Number of Rounds [BG = 500,Dynamic BG Cover[45-55]%,Binomial mix]

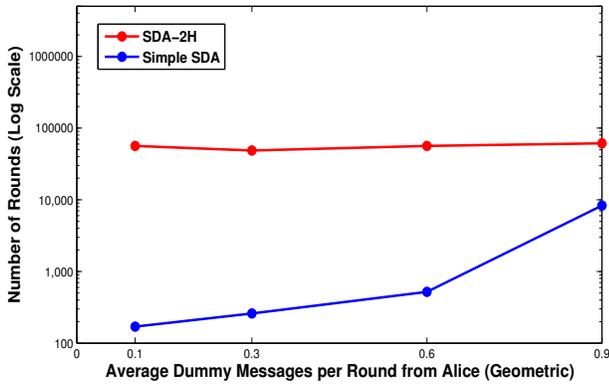


Figure 7: Number of Rounds [BG = 500,Non-dynamic BG Cover[100%],Binomial mix]

and actually distracts from SDA’s effectiveness. In practice, the attacker would use whichever attack was most effective for his estimate of the system parameters.

In reality, SDA-2H may not be more powerful than SDA in all the cases, but SDA-2H is an effective attack model to explore the effect of background cover, as it successfully reduces the effect of Alice Cover.

6. RELATED WORK

Significant work was made in this area in 2004, when Mathewson and Dingledine extend previous work on SDA to show that even in a binomial mix with a non-global adversary, intersection attacks can still be successful, but may be made to require an extremely long period of observation [9]. They accomplish this in a manner similar to ours: by varying the amount of background traffic and padding (cover traffic). Their findings confirm that to achieve the greatest anonymity, the delay ($pDelay$ in our models) must be highly varied, padding must be consistent, and the attacker must have a limited view of the network, with the inability to distinguish between times when the victim is active and not (*i.e.*, the victim rarely, if ever, goes offline, and consistently sends cover traffic).

Danezis et. al. have presented a model, TS-SDA “...that tries to uncover the receivers of messages sent through an anonymizing network supporting anonymous replies. [*e.g.*, Mix minion [5]]”

[4]. The work proposes a linear approximation technique for predicting likely receivers, and demonstrates that TS-SDA’s effectiveness exceeds that of SDA, since it considers the replies as well, hence a “two sided” attack. The linear approximation uses a complex weighting to decide the relationship of a message to Alice. (*i.e.*, did she initiate the conversation, or was she just replying to a larger group message?) Our technique would likely not apply, as it depends on a simple derivation of the total number of messages sent by the mix, which becomes a complex problem when senders double as receivers.

A much more recent initiative into SDA has been undertaken by Mallesh and Wright by viewing anonymous communication over the mix as two-way instead of the traditional one-way approach of SDA [8]. This attack, referred to by the authors as the reverse statistical disclosure attack (RSDA), is unique in its approach. The attacker basically applies the SDA attack to all users in a mix, thus observing which users Alice is sending messages to and which users are sending messages to Alice. By observing these two sets of information, the attacker can obtain a profile of Alice and locate contacts of Alice that may be missed by traditional SDA. This attack proves to be much more reliable and faster than SDA. Readers may note that Danezis et. al. [4] employ an attack which seems to have a similar “two sided” attack structure. However, the attack presented by Danezis differs from RSDA in that it is only interested in receivers to whom Alice initiates the message, while RSDA is interested in any contact of Alice, regardless of who initiated the conversation.

The innovative approaches offered by RSDA and TS-SDA unfortunately can not be applied to our currently proposed attack model as SDA-2H shares a glaring weakness of SDA, in that it can not distinguish between replies and initiated messages (*i.e.*, it is one-sided). While the models of Danezis and Mallesh are noteworthy, the increasing importance of peer-to-peer communication does not imply that improvements for cases where anonymity is still one-sided are useless; but rather, expands the applicability of mix networks in new directions.

7. FUTURE WORK

Introducing error into the components of the SDA equation has shed light on the anonymity of mixes, though the exact effects have not been explored in a mathematically rigorous way. By returning to the original formulation of [3], we will analyze the effect on the lower bound given for the number of rounds an attacker can expect to take for a varying amount of error in the number of messages. Namely, Danezis gives this lower bound, t (with a given “security parameter”, l), as,

$$t > \left[m \cdot l \left(\sqrt{\frac{N-1}{N} (b-1)} + \sqrt{\frac{N-1}{N^2} (b-1) + \frac{m-1}{m}} \right) \right]^2 \quad (4)$$

he also derives a precondition for an attacker’s success based on Alice’s sending volume, $m < \frac{N}{b-1}$. Our continued interest is in understanding the precise effect on these bounds for a given amount of padding, both by Alice and the background; as the equations alone seem not to incorporate the additional information an approach such as SDA-2H provides.

One noted weakness of SDA is that it relies on a relatively Alice-free background against which to make its correlations. Danezis notes, however, that so long as Alice’s behavior varies substantially, even a continuous sending pattern does not provide complete anonymity; it would have to abide by certain conditions – which we intend to compare with SDA-2H by varying the three types of cover discussed in section 2. We expect that as we approach an

equal amount of inbound and outbound messages, we may come closer to providing better anonymity for all the senders, not just ones that send sufficiently more than their peers.

Using this method, $pDelay$ is not considered, so it is unlikely to improve anonymity as much as its delayed users may hope; though we suspect it will fare much better when combined with receiver-bound cover.

8. CONCLUSIONS

In this paper, we have described and evaluated SDA-2H. We specifically use SDA-2H as a tool to measure the previously unknown effects of background cover on the anonymity of mixed based systems. By using a carefully chosen set of values for our attack simulator (to guarantee maximum coverage of data values) we have finally understood some potential benefits of adding background cover to a mix based system.

We also show how the same background cover traffic can be nullified by a well-informed attacker. Using the information gleaned from these experiments, coupled together with a greater understanding of mixes, we can be one step closer to obtaining the ideal form of anonymous communication, one that is insusceptible to any attack.

9. REFERENCES

[1] D. L. Chaum. Untraceable electronic mail, return addresses, and digital pseudonyms. *Communications of the ACM*, 24, February 1981.

[2] G. Danezis. Statistical disclosure attacks: Traffic confirmation in open environments. In *Proc. Security and Privacy in the Age of Uncertainty (Sec2003)*, May 2003.

[3] G. Danezis. Better Anonymous Communications. PhD thesis, University of Cambridge, July 2004.

[4] G. Danezis, C. Diaz, and C. Troncoso. Two-sided statistical disclosure attack. In *Proc. Privacy Enhancing Technologies (PET '07)*, June 2007.

[5] G. Danezis, R. Dingledine, and N. Mathewson. Mixminion: Design of a type iii anonymous remailer protocol. In *Proc. IEEE Symposium on Security and Privacy (Oakland '03)*, May 2003.

[6] D. Kesdogan, D. Agrawal, and S. Penz. Limits of anonymity in open environments. In *Proc. Information Hiding (IH '02)*, October 2002.

[7] N. Mallesh and M. Wright. Countering statistical disclosure with receiver-bound cover traffic. In *Proc. European Symposium on Research in Computer Security (ESORICS '07)*, September 2007.

[8] N. Mallesh and M. Wright. The reverse statistical disclosure attack. In *Proc. Information Hiding (IH '10)*, June 2010.

[9] N. Mathewson and R. Dingledine. Practical traffic analysis: Extending and resisting statistical disclosure. In *Proc. Privacy Enhancing Technologies (PET '04)*, May 2004.

[10] P. Venkatasubramanian and V. Anantharam. On the anonymity of chaum mixes. In *Proc. International Symposium on Information Theory*, July 2008.