

Five Practical Attacks for “Optimistic Mixing for Exit-Polls”

Douglas Wikström

Swedish Institute of Computer Science (SICS)
Box 1263, S-164 29 Kista, Sweden
`douglas@sics.se`

Abstract. Golle, Zhong, Boneh, Jakobsson, and Juels [9] recently presented an efficient mix-net, which they claim to be both robust and secure. We present five practical attacks for their mix-net, and break both its privacy and robustness.

The first attack breaks the privacy of any given sender *without corrupting any mix-server*. The second attack requires that the first mix-server is corrupted. Both attacks are adaptations of the “relation attack” introduced by Pfitzmann [24, 23].

The third attack is similar to the attack of Desmedt and Kurosawa [4] and breaks the privacy of *all* senders. It requires that all senders are honest and that the last mix-server is corrupted.

The fourth attack may be viewed as a novel combination of the ideas of Lim and Lee [16] and Pfitzmann [24, 23]. It breaks the privacy of any given sender, and requires that the first and last mix-servers are corrupted. This attack breaks also Jakobsson [14], including the fixed version of Mitomo and Kurosawa [18].

The fifth attack breaks the robustness in a novel way. It requires corruption of some senders and the first mix-server. This attack breaks also Jakobsson and Juels [15].

1 Introduction

The notion of a mix-net was invented by Chaum [3] and further developed by a number of people. Properly constructed a mix-net enables a set of senders to send messages anonymously.

Informally the requirements on a mix-net are: correctness, privacy, robustness, availability, and efficiency. Correctness implies that the result is correct given that all mix-servers are honest. Privacy implies that if a fixed minimum number of mix-servers are honest privacy of the sender of a message is ensured. Robustness implies that if a fixed number of mix-servers are honest, then any attempt to cheat is detected and defeated.

A mix-net consists of a number of mix-servers that collectively execute a protocol. The idea is that each mix-server receives a list of encrypted messages, transforms them, using partial decryption or random re-encryption, reorders them, and then outputs the randomly transformed and reordered list of cleartext messages. The secret permutation is shared among the mix-servers.

1.1 Previous Work and Applications of Mix-Nets

The mixing paradigm has been used to accomplish privacy in many different scenarios. Chaum’s original “anonymous channel” [3, 21] enables a sender to securely send mail anonymously. When constructing election schemes [3, 6, 22, 25, 20] the mix-net is used to ensure that the vote of a given voter can not be revealed. Also in the construction of electronic cash systems [13] mix-nets have been used to ensure privacy.

Abe gives an efficient construction of a general mix-net [1] and argues about its properties. Jakobsson has written a number of more general papers on the topic of mixing [12, 14] also focusing on efficiency.

There has been a breakthrough in the construction of zero-knowledge proofs of a correct shuffle recently. Furukawa and Sako [7], and Neff [19] respectively have both found efficient ways to design such proofs.

A new approach to practical mixing is given by Golle et al. [9]. They combine a robustness test based partly on work by Jakobsson and Juels, with the notion of “double enveloping”. The latter notion is introduced independently by Wikström [29], except that he uses different keys for the two layers, and a proof of knowledge of the inner most cleartext.

Desmedt and Kurosawa [4] describe an attack on a protocol by Jakobsson [12]. Similarly Mitomo and Kurosawa [18] exhibit a weakness in another protocol by Jakobsson [14]. Pfitzmann has given some general attacks on mix-nets [24, 23], and Michels and Horster give additional attacks in [17].

This paper is based on two technical reports [30, 31]. Apparently Abe [2] has independently found attacks similar to those of the first of our technical reports, i.e. the first two attacks in this paper.

2 Review of “Optimistic Mixing for Exit-Polls”

We present a short review of the relevant parts of the protocol of Golle et al. [9]. The description given here is as close as possible to the original, but we avoid details irrelevant to our attacks and change some notation to simplify the exposition of the attacks. For details we refer the reader to [9].

2.1 Participants and Setup

The protocol assumes the existence of a bulletin board on which each participant has a dedicated area on which only she can write. No participant can erase anything on the bulletin board, but all participants can read everything.

The participants of the protocol are N senders, and a relatively small number of *mix-servers*, M_1, \dots, M_k . Each sender encrypts its message, and writes it on the bulletin board. The mix-servers then execute the mix-net protocol.

The protocol employs an El Gamal [5] cryptosystem in a subgroup G_Q of prime order Q of the multiplicative group modulo a prime P , i.e. \mathbb{Z}_P^* . A private key x is generated by choosing $x \in \mathbb{Z}_Q$ uniformly and independently at

random. The corresponding public key is (g, y) , where g is a generator of G_Q , and $y = g^x$. Encryption of a message $m \in G_Q$ using the public key (g, y) is given by $E_{(g,y)}(m, r) = (g^r, y^r m)$, where r is chosen uniformly at random from \mathbb{Z}_Q , and decryption of a cryptotext on the form $(u, v) = (g^r, y^r m)$ using the private key x is given by $D_x(u, v) = u^{-x} v = m$. The El Gamal system also has the re-encryptability¹ property, i.e. given (u, v) and the public key (g, y) , anybody can “update” the randomness used in the encryption of m , by forming $(g^{r'} u, y^{r'} v) = (g^{r+r'}, y^{r+r'} m) = E_{(g,y)}(m, r + r')$.

In the setup stage each mix-server M_j is somehow given a random $x_j \in \mathbb{Z}_Q$, and $y_l = g^{x_l}$ for $l \neq j$. The value x_j is also shared with the other mix-servers using a threshold verifiable secret sharing (VSS) scheme². Thus, if any mix-server M_j is deemed to be cheating the other mix-servers can verifiably reconstruct its private key x_j . The mix-servers can also compute $y = \prod_{j=1}^k y_j$, which gives a joint public key (g, y) , with secret corresponding private key $x = \sum_{j=1}^k x_j$.

A practical advantage of the mix-net is that it can be used to execute several mix-sessions using the same set of keys, i.e. the El Gamal keys are not changed between mix-sessions. To be able to do this the proofs of knowledge below are bound to a mix-session identifier id that is unique to the current mix-session.

2.2 Sending a Message to the Mix-Net

A typical honest sender, Alice, computes the following to send a message m to the mix-net:

$$(u, v) = E_{(g,y)}(m), \quad w = h(u, v), \quad \text{and} \\ \alpha = [E_{(g,y)}(u), E_{(g,y)}(v), E_{(g,y)}(w)] = [(\mu_1, \mu_2), (\nu_1, \nu_2), (\omega_1, \omega_2)] ,$$

where $h : \{0, 1\}^* \rightarrow G_Q$ is a hash function modeled by a random oracle. Then Alice computes a zero-knowledge proof of knowledge $\pi_{\text{id}}(u, v, w)$, in the random oracle model of u, v and w , which depends on the current mix-session identifier id. Finally Alice writes $(\alpha, \pi_{\text{id}}(u, v, w))$ on the bulletin board.

2.3 Execution of the Mix-Net

First the mix-servers remove any duplicate inputs to the mix-net, and discard input tuples that contain components not in the subgroup G_Q . Then the mix-servers discard all input tuples where the proof of knowledge is not valid for the current mix-session. Let $L_0 = \{[(a_{0,i}, b_{0,i}), (c_{0,i}, d_{0,i}), (e_{0,i}, f_{0,i})]\}_{i=1}^N$ be the resulting list of triples of El Gamal pairs. The mixing then proceeds in the following stages.

¹ Related terms are “homomorphic property”, or “malleability property”.

² Golle et al. [9] discuss different variants for sharing keys, but we choose to present a simple variant, since it has no impact on our attacks.

First Stage: Re-randomization and Mixing. This step proceeds as in all re-randomization mix-nets based on El Gamal. One by one, the mix-servers M_1, \dots, M_k randomize all the inputs and their order. (Note that the components of triples are not separated from each other during the re-randomization.) In addition, each mix-net must give a proof that the product of the plaintexts of all its inputs equals the product of the plaintexts of all its outputs. The protocol proceeds as follows.

1. Each mix-server M_j reads from the bulletin board the list L_{j-1} output by the previous mix-server.
2. The mix-server then chooses $r_{ji}, s_{ji}, t_{ji} \in \mathbb{Z}_Q$, for $i = 1, \dots, N$, randomly and computes the re-randomized list:

$$\{[(g^{r_{ji}} a_{j-1,i}, y^{r_{ji}} b_{j-1,i}), (g^{s_{ji}} c_{j-1,i}, y^{s_{ji}} d_{j-1,i}), (g^{t_{ji}} e_{j-1,i}, y^{t_{ji}} f_{j-1,i})]\}_{i=1}^N .$$

The above list of triples is then randomly permuted, and the resulting list: $L_j = \{(a_{j,i}, b_{j,i}), (c_{j,i}, d_{j,i}), (e_{j,i}, f_{j,i})\}_{i=1}^N$ is written on the bulletin board.

3. Define $a_j = \prod_{i=1}^N a_{j,i}$, and define b_j, c_j, d_j, e_j , and f_j correspondingly. The mix-server proves in zero-knowledge that $\log_g a_j/a_{j-1} = \log_y b_j/b_{j-1}$, $\log_g c_j/c_{j-1} = \log_y d_j/d_{j-1}$, and $\log_g e_j/e_{j-1} = \log_y f_j/f_{j-1}$. This implies that $D_x(a_j, b_j) = D_x(a_{j-1}, b_{j-1})$, and similarly for the pairs (c_j, d_j) and (e_j, f_j) , i.e. the component-wise product of the inner triples remains unchanged by the mix-server.

Remark 1. Since $\log_y b_j/b_{j-1} = \log_g a_j/a_{j-1} = \sum_{i=1}^N r_{ji}$, and M_j knows the latter sum, the proof in Step 3) can be implemented by a standard zero-knowledge proof of knowledge in the random oracle model, and similarly for the pairs (c_j, d_j) , and (e_j, f_j) . Such proofs are similar to Schnorr signatures [26].

Second Stage: Decryption of the Inputs.

1. A quorum of mix-servers jointly decrypt each triple of ciphertexts in L_k to produce a list L on the form $L = \{(u_i, v_i, w_i)\}_{i=1}^N$. Since the method used to do this is irrelevant to our attacks, we do not present it here.
2. All triples for which $w_i = h(u_i, v_i)$ are called *valid*.
3. Invalid triples are investigated according to the procedure described below. If the investigation proves that all invalid triples are *benign* (only senders cheated), we proceed to Step 4. Otherwise, the decryption is aborted, and we continue with the back-up mixing.
4. A quorum of mix-servers jointly decrypt the ciphertexts (u_i, v_i) for all valid triples. This successfully concludes the mixing. The final output is defined as the set of plaintexts corresponding to valid triples.

Special Step: Investigation of Invalid Triples. The mix-servers must reveal the path of each invalid triple through the various permutations. For each

invalid triple, starting from the last server, each server reveals which of its inputs corresponds to this triple, and how it re-randomized this triple. One of two things may happen:

- **Benign case (only senders cheated):** if the mix-servers successfully produce all such paths, the invalid triples are known to have been submitted by users. The decryption resumes after the invalid triples have been discarded.
- **Serious case (one or more servers cheated):** if one or more servers fail to recreate the paths of invalid triples, these mix-servers are accused of cheating and replaced, and the mix-net terminates producing no output. In this case, the inputs are handed over to the back-up mixing procedure below.

Back-Up Mixing. The *outer-layer* encryption of the inputs posted to the mix-net is decrypted by a quorum of mix-servers. The resulting list of *inner-layer* ciphertexts becomes the input to a standard re-encryption mix-net based on El Gamal (using, for example Neff’s scheme described in [19]). Then the output of the standard mix-net is given as output by the mix-net.

Remark 2. It is impossible to find two lists $\{(u_i, v_i)\}_{i=1}^N \neq \{(u'_i, v'_i)\}_{i=1}^N$ such that $\prod_{i=1}^N h(u_i, v_i) = \prod_{i=1}^N h(u'_i, v'_i)$, if the product is interpreted in a group where the discrete logarithm problem is hard. This is stated as a theorem by Wagner³ [28], and appears as a lemma in Golle et al. [9].

During the re-encryption and mixing stage each mix-server proves in zero-knowledge that it leaves the component-wise product $(\prod u_i, \prod v_i, \prod w_i)$, of the inner triples (u_i, v_i, w_i) unchanged, but individual triples may still be corrupted. Then invalid triples are traced back. This leaves only valid inner triples in the output and the proofs of knowledge of each server are used to conclude that the component-wise product of these valid inner triples was left unchanged by the mix-net. Golle et al. [9] then refer to the lemma and conclude that the set of valid triples in the output is identical to the set of valid triples hidden in the double encrypted input to the mix-net.

Unfortunately, this intuitively appealing construction is flawed as we explain in Section 3.5. Furthermore, our third attack in Section 3.3 shows that it is possible to cheat without changing the set of inner triples.

3 The Attacks

The goal of the adversary Eve is to break the privacy of our typical honest sender Alice and cheat without detection. Each of our attacks illustrates a separate weakness of the protocol. The first two attacks are adaptations of the “relation attack”, introduced by Pfitzmann [24, 23], to the setting with double enveloping. The idea of the “relation attack” is that to break the privacy of Alice, Eve computes a cryptotext of a message related to Alice’s message. Then the mix-net is run as usual. The output of the mix-net contains two messages related in

³ Wagner credits Wei Dai with this observation.

a way chosen by Eve. Some relations enable Eve to determine the message sent by Alice. The third attack is similar to the attack of Desmedt and Kurosawa [4] in that it exploits intermediate results of the protocol and fools a “product test”. The fourth attack may be viewed as a novel combination of the ideas of Lim and Lee [16] and Pfitzmann [24, 23]. The fifth attack is novel.

3.1 First Attack: Honest Mix-Servers

We show that the adversary Eve can break the privacy of the typical sender Alice. All that is required is that Eve can send two messages to the mix-net, which is a natural assumption in most scenarios. In the following we use the notation for the cryptotext of Alice introduced in Section 2.2. Eve does the following:

1. Eve chooses δ and γ randomly in \mathbb{Z}_Q , and computes:

$$w_\delta = h(\mu_1^\delta, \mu_2^\delta), \quad \alpha_\delta = (E_{(g,y)}(\mu_1^\delta), E_{(g,y)}(\mu_2^\delta), E_{(g,y)}(w_\delta)) \quad , \quad \text{and}$$

$$w_\gamma = h(\mu_1^\gamma, \mu_2^\gamma), \quad \alpha_\gamma = (E_{(g,y)}(\mu_1^\gamma), E_{(g,y)}(\mu_2^\gamma), E_{(g,y)}(w_\gamma)) \quad .$$

Then Eve computes corresponding proofs of knowledge $\pi_{\text{id}}(\mu_1^\delta, \mu_2^\delta, w_\delta)$ and $\pi_{\text{id}}(\mu_1^\gamma, \mu_2^\gamma, w_\gamma)$. This gives Eve two perfectly valid pairs $(\alpha_\delta, \pi_{\text{id}}(\mu_1^\delta, \mu_2^\delta, w_\delta))$, $(\alpha_\gamma, \pi_{\text{id}}(\mu_1^\gamma, \mu_2^\gamma, w_\gamma))$, that she sends to the bulletin board (possibly by corrupting two senders).

2. Eve waits until the mix-net has successfully completed its execution. During the execution of the mix-net the mix-servers first jointly decrypt the “outer layer” of the double encrypted messages. After benign tuples have been removed the result is a list of valid triples

$$((u_1, v_1, w_1), \dots, (u_N, v_N, w_N)) \quad . \tag{1}$$

The final output of the mix-net is the result of decrypting each inner El Gamal pair (u_i, v_i) and results in a list of cleartext messages (m_1, \dots, m_N) .

3. Eve computes the list $(m'_1, \dots, m'_N) = (m_1^{\delta/\gamma}, \dots, m_N^{\delta/\gamma})$, and then finds a pair (i, j) such that $m_i = m'_j$. From this she concludes that with very high probability $m_j = u^\gamma$. Then she computes $z = m_j^{1/\gamma}$, and finds a triple (u_l, v_l, w_l) in the list (1) such that $z = u_l$. Finally she concludes that with very high probability m_l was the message sent by Alice to the mix-net.

Remark 3. At additional computational cost it suffices for Eve to send 2 messages to break the privacy of K senders. Suppose Eve wants to break the privacy also of Bob who sent m' encrypted as $(u', v') = E_{(g,y)}(m')$, $w' = h(u', v')$, and $\alpha' = [E_{(g,y)}(u'), E_{(g,y)}(v'), E_{(g,y)}(w')] = [(\mu'_1, \mu'_2), (\nu'_1, \nu'_2), (\omega'_1, \omega'_2)]$. Then Eve performs the attack above with the change that she starts with a single pair $(\mu_1^\zeta, \mu_2^\zeta)$ for some randomly chosen ζ instead of the two distinct pairs (μ_1, μ_2) , and (μ'_1, μ'_2) that would have given two “unrelated” attacks. The original third step of the attack first gives Eve $z = u^\zeta u'$. To finish the attack she finds a pair (l, l') such that $u_l^\zeta u_{l'} = z$, and concludes that with high probability Alice sent m_l and Bob sent $m_{l'}$. The approach is generalized to higher dimensions in the natural way to break the privacy of several senders (K must clearly be reasonably sized).

Why the Attack Is Possible. The attack exploits two different flaws of the protocol. The first is that the sender of a message, e.g. Alice, proves only knowledge of the inner El Gamal pair (u, v) and the hash value $w = h(u, v)$, and not knowledge of the message m . The second flaw is that identical El Gamal keys are used for both the inner and outer El Gamal system.

Anybody can compute a single encrypted message $(\mu_1^\delta, \mu_2^\delta) = (g^{r\delta}, y^{r\delta}u^\delta) = E_{(g,y)}(u^\delta, r\delta)$ of a power u^δ of the first component u of the *inner* El Gamal pair (u, v) of the triple α sent by Alice. Anybody can also compute a proof of knowledge of $(\mu_1^\delta, \mu_2^\delta)$ and $w_\delta = h(\mu_1^\delta, \mu_2^\delta)$ (and similarly for $(\mu_1^\gamma, \mu_2^\gamma)$ and w_γ).

The first flaw allows Eve to input triples of El Gamal pairs with such proofs of knowledge to the mix-net. The second flaw allows Eve to use the mix-net to decrypt $(\mu_1^\delta, \mu_2^\delta)$, and thus get her hands on u^δ (and similarly for u^γ). Eve can then identify (u, v) as the inner El Gamal pair of Alice and break her privacy.

3.2 Second Attack: Different Keys and Corrupt Mix-Server

Suppose we change the protocol slightly by requiring that the mix-servers generate separate keys for the outer and inner El Gamal systems, to avoid the first attack of Section 3.1. We assume that there are two different key pairs $((g, y_{\text{in}}), x_{\text{in}})$ and $((g, y_{\text{out}}), x_{\text{out}})$, for the inner and outer El Gamal layers respectively. We also assume that these keys have been shared similarly as the original key pair $((g, y), x)$. This is the type of double enveloping proposed by Wikström [29]. For the second attack to succeed we need some additional assumptions.

Unclear Details and Additional Assumptions. We start by quoting Section 5, under “Setup.” point 4 of Golle et al. [9], which presents the proof of knowledge $\pi_{\text{id}}(u, v, w)$ of the sender Alice:

4. This proof of knowledge should be bound to a unique mix-session identifier to achieve security over multiple invocations of the mix. Any user who fails to give the proof is disqualified, and the corresponding input is discarded.

If different keys are used for each mix-session, then the above makes no sense, since the proof of knowledge of u, v and w already depends on the public key of the outer El Gamal system. There is clearly practical value in not changing keys between mix-sessions. We assume that the keys are not changed between mix-sessions even if a mix-server is found to be cheating. If a mix-server is found to be cheating, its shared keys are instead reconstructed by the remaining mix-servers using the VSS-scheme, and in later mix-sessions the actions of the cheating mix-server are performed in the open (the details of this does not matter to our attack). Under these assumptions we can give an attack on the protocol.

The original paper of Golle et al. [9] does not explicitly say if the discovery of the corrupted mix-server results in a new execution of the key generation

protocol. Apparently the intention of the authors is to let the remaining mix-servers generate a new set of keys if any cheating is discovered [10].

The attack is interesting even though this interpretation is not the one intended by the authors, since it shows the importance of explicitly defining all details of protocols and highlights some issues with running several concurrent mix-sessions using the same set of keys.

The Attack. Apart from the above assumptions, the attack only requires that the first mix-server in the mix-chain is corrupted. The attack is employed during two mix-sessions using the same keys and the corrupted mix-server is identified as a cheater in the first mix-session. In the following we describe the actions of Eve during the first and second mix-sessions, respectively.

THE FIRST MIX-SESSION. We assume that Alice and some other arbitrary sender Bob have sent inputs to the mix-net (and use the notation of Remark 3 for the input of Bob). Eve corrupts M_1 . It then replaces α and α' with: $[E_{y_{\text{out}}}(u), E_{y_{\text{out}}}(v), E_{y_{\text{out}}}(w')]$, and $[E_{y_{\text{out}}}(u'), E_{y_{\text{out}}}(v'), E_{y_{\text{out}}}(w)]$ respectively, in its input list, i.e. the third components of the two triples are shifted. Then Eve forces M_1 to simulate a completely honest mix-server on the resulting altered list $L'_0 = \{[(a'_{0,i}, b'_{0,i}), (c'_{0,i}, d'_{0,i}), (e'_{0,i}, f'_{0,i})]\}_{i=1}^N$. Note that $\prod_{i=1}^N a'_{0,i} = a_0$, and similarly for b_0, c_0, d_0, e_0 , and f_0 . Thus the simulated honest mix-server outputs perfectly valid zero-knowledge proofs that the product of the inner triples are unchanged.

At the end of the mixing the mix-servers verify the tuples and discover the invalid tuples (u, v, w') and (u', v', w) . These tuples are traced back all the way to the first mix-server, which is revealed as a cheater. In this process Eve is able to link Alice to (u, v) (and Bob to (u', v')). Finally the honest mix-servers finish the protocol by using the general constructions based on the work by Neff [19] as in Golle et al. [9].

THE SECOND MIX-SESSION. To allow the mix-net to execute a second mix-session using the same set of keys, the cheater’s key is reconstructed and revealed by a quorum of the mix-servers.

To determine the contents of the El Gamal pair (u, v) , Eve waits for the second mix-session using the same set of keys. Then she uses a “relation attack” [24, 23, 12] in the second mix-session to decrypt (u, v) , i.e. Eve does the following:

1. Eve chooses δ and γ randomly in \mathbb{Z}_Q , and computes:
 $w_\delta = h(u^\delta, v^\delta)$, $\alpha_\delta = (E_{y_{\text{out}}}(u^\delta), E_{y_{\text{out}}}(v^\delta), E_{y_{\text{out}}}(w_\delta))$, and
 $w_\gamma = h(u^\gamma, v^\gamma)$, $\alpha_\gamma = (E_{y_{\text{out}}}(u^\gamma), E_{y_{\text{out}}}(v^\gamma), E_{y_{\text{out}}}(w_\gamma))$.
 Then Eve computes corresponding proofs of knowledge $\pi_{\text{id}}(u^\delta, v^\delta, w_\delta)$ and $\pi_{\text{id}}(u^\gamma, v^\gamma, w_\gamma)$. This gives Eve two perfectly valid pairs $(\alpha_\delta, \pi_{\text{id}}(u^\delta, v^\delta, w_\delta))$, $(\alpha_\gamma, \pi_{\text{id}}(u^\gamma, v^\gamma, w_\gamma))$, which she sends to the bulletin board (possibly by corrupting two senders).
2. Eve waits until the mix-net has successfully completed its execution. The final output of the mix-net is a list of cleartext messages (m_1, \dots, m_N) .

3. Note that $m_i = m^\delta$ and $m_j = m^\gamma$ for some i and j . Eve computes $\delta\gamma^{-1} \bmod Q$, computes the list $(m'_1, \dots, m'_N) = (m_1^{\delta/\gamma}, \dots, m_N^{\delta/\gamma})$, and finally finds a pair (i, j) such that $m_i = m'_j$. Then she concludes that with high probability $m_j^{1/\gamma}$ is the message sent by Alice to the mix-net in the *first* mix-session.

Remark 4. The attack is easily generalized to break the privacy of several senders by using a circular shift of the third components during the first mix-session. It is also easy to see that Remark 3 can be applied to reduce the number of messages sent by Eve during the second mix-session.

Why the Attack Is Possible. The attack exploits that the sender of a message only proves knowledge of the inner triple (u, v, w) . At the cost of detected cheating Eve finds a (u, v) corresponding to Alice, and then uses the second mix-session as a decryption oracle to get her hands on m .

A Note on Concurrent Mix-Sessions. Ignoring the other attacks, a simple countermeasure to the second attack above, is to stipulate that if a cheating mix-server is identified new keys must be generated for the next mix-session.

A disadvantage of this countermeasure is that the mix-net can not be allowed to execute several concurrent mix-sessions using the same keys. If one mix-session is still receiving messages while another mix-session discovers a cheating mix-server the second attack of Section 3.2 can still be applied. The problem is not solved by running the back-up mix-net of Neff [19] on all mix-sessions using the same keys at this point.

This problem of concurrency may seem academic, since in most election scenarios it is not very cumbersome to have different keys for each mix-session, but in future applications of mix-nets it may be useful to run several concurrent mix-sessions using the same keys.

3.3 Third Attack: Honest Senders and One Corrupt Mix-Servers

In this section we assume that all senders and all mix-servers, except the last mix-server M_k , are honest. The last mix-server M_k is corrupted by the adversary Eve and performs the attack. The attack breaks both robustness and privacy.

To simplify our notation we write $L_0 = \{\alpha_i\}_{i=1}^N$ for the input list, where we define $\alpha_i = [(a_{0,i}, b_{0,i}), (c_{0,i}, d_{0,i}), (e_{0,i}, f_{0,i})]$ to be the tuple sent by sender S_i . Instead of following the protocol, M_k proceeds as follows in the first stage.

1. It computes the 6-tuple: $(a', b', \dots, f') = (a_{k-1}/a_0, b_{k-1}/b_0, \dots, f_{k-1}/f_0)$, and the tuple $\alpha'_1 = [(a'a_{0,1}, b'b_{0,1}), (c'c_{0,1}, d'd_{0,1}), (e'e_{0,1}, f'f_{0,1})]$.
2. Then it forms the list $L'_{k-1} = \{\alpha'_1, \alpha_2, \dots, \alpha_N\}$, i.e. a copy of L_0 with the first tuple α_1 replaced by α'_1 .
3. When M_k is supposed to re-randomize and permute the output L_{k-1} of M_{k-1} it instead simulates the actions of an honest mix-server on the corrupted input L'_{k-1} . The output list written to the bulletin board by the simulated mix-server is denoted L_k .

4. Eve waits until the inner layer has been decrypted and uses her knowledge of the permutation that relates L_k to L_0 to break the privacy of all senders.

We show that the attack goes undetected, i.e. the mix-servers decrypt the inner layer. This implies that the attack succeeds.

Firstly, consider the proof of knowledge that M_k produces during the re-encryption and mixing stage. Define $a'_{k-1} = (a'a_{0,1}) \prod_{i=2}^N a_{0,i}$, and similarly for b'_{k-1} , c'_{k-1} , d'_{k-1} , e'_{k-1} , and f'_{k-1} . In Step 3 above, the simulated honest mix-server outputs proofs of knowledge of the logarithms: $\log_g a_k/a'_{k-1} = \log_y b_k/b'_{k-1}$, $\log_g c_k/c'_{k-1} = \log_y d_k/d'_{k-1}$, and $\log_g e_k/e'_{k-1} = \log_y f_k/f'_{k-1}$. By construction we have that $a'_{k-1} = (a'a_{0,1}) \prod_{i=2}^N a_{0,i} = a' \prod_{i=1}^N a_{0,i} = \frac{a_{k-1}}{a_0} a_0 = a_{k-1}$, and similarly for b_{k-1} , c_{k-1} , d_{k-1} , e_{k-1} , and f_{k-1} . This implies that the proof of knowledge produced by M_k is deemed valid by the other mix-servers.

Secondly, consider the investigation of invalid inner triples. Tracing back invalid triples is difficult to M_k , since it does not know the re-encryption exponents and the permutation relating L_{k-1} and L_k . We show that there are no invalid inner triples. Suppose we define the sums $r = \sum_{j=1}^{k-1} \sum_{i=1}^N r_{ji}$, $s = \sum_{j=1}^{k-1} \sum_{i=1}^N s_{ji}$, and $t = \sum_{j=1}^{k-1} \sum_{i=1}^N t_{ji}$, i.e. we form the sum of all re-encryption exponents used by all mix-servers except the last, for the first second and third El Gamal pairs respectively. Since all mix-servers except M_k are honest, we have $(a', b', c', d', e', f') = (g^r, y^r, g^s, y^s, g^t, y^t)$, which implies that α'_1 is a valid re-encryption of α_1 . Thus M_k does not corrupt any inner triple by simulating an honest mix-server on the input L'_{k-1} . Since all senders are honest and the set of inner triples hidden in L_0 and L'_{k-1} are identical, there are no invalid inner triples. Thus cheating is not detected and robustness is broken.

We conclude that the mix-servers decrypt the inner triples, i.e. the privacy of *all* senders is broken.

Why the Attack Is Possible. The third attack above exploits that the last mix-server M_k is not forced to take the output L_{k-1} of the next to last mix-server as input. This allows M_k to use a slightly modified version of L_0 instead, which breaks the privacy of all senders.

3.4 Fourth Attack: Two Corrupt Mix-Servers

In this section we assume that the first and last mix-servers, M_1 and M_k , are corrupted. We give a novel attack that breaks the privacy of any given sender Alice. Let $L_0 = \{\alpha_i\}_{i=1}^N$, where $\alpha_i = [(a_{0,i}, b_{0,i}), (c_{0,i}, d_{0,i}), (e_{0,i}, f_{0,i})]$. Without loss we let α_1 and α_2 be the tuples sent by Alice and Bob respectively. Let $\xi \in \mathbb{Z}_P^* \setminus G_Q$. Eve corrupts M_1 and M_k , and they proceed as follows.

1. M_1 computes the elements $\alpha'_1 = [(\xi a_{0,1}, b_{0,1}), (c_{0,1}, d_{0,1}), (e_{0,1}, f_{0,1})]$ and $\alpha'_2 = [(\xi^{-1} a_{0,2}, b_{0,2}), (c_{0,2}, d_{0,2}), (e_{0,2}, f_{0,2})]$, and forms the modified list $L'_0 = \{\alpha'_1, \alpha'_2, \alpha_3, \dots, \alpha_N\}$. Then M_1 simulates an honest mix-server on input L'_0 .

2. M_k simulates an honest mix-server on input L_{k-1} , but it does not write the output L_k of this simulation on the bulletin board. Let $L_k = \{\beta_i\}_{i=1}^N$, where $\beta_i = [(a_{k,i}, b_{k,i}), (c_{k,i}, d_{k,i}), (e_{k,i}, f_{k,i})]$. M_k finds $l, l' \in \{1, \dots, N\}$ such that $a_{k,l}^Q = \xi^Q$ and $a_{k,l'}^Q = \xi^{-Q}$.
Then it computes the tuples $\beta'_l = [(\xi^{-1}a_{k,l}, b_{k,l}), (c_{k,l}, d_{k,l}), (e_{k,l}, f_{k,l})]$ and $\beta'_{l'} = [(\xi a_{k,l'}, b_{k,l'}), (c_{k,l'}, d_{k,l'}), (e_{k,l'}, f_{k,l'})]$, and writes L'_k , where $L'_k = \{\beta_1, \dots, \beta_{l-1}, \beta'_l, \beta_{l+1}, \dots, \beta_{l'-1}, \beta'_{l'}, \beta_{l'+1}, \dots, \beta_N\}$, on the bulletin board.
3. The mix-net outputs (m_1, \dots, m_N) and Eve concludes that Alice sent m_l .

Since all mix-servers except M_1 and M_k are honest there exists $l, l' \in \{1, \dots, N\}$ and $r, r' \in \mathbb{Z}_Q$ such that $a_{k,l} = g^r \xi a_{0,1}$ and $a_{k,l'} = g^{r'} \xi^{-1} a_{0,2}$. This implies that $a_{k,l}^Q = \xi^Q (g^r a_{0,1})^Q = \xi^Q$ and $a_{k,l'}^{-Q} = \xi^{-Q}$, since $\beta^Q = 1$ for any $\beta \in G_Q$. We have $\xi^Q \neq 1 \neq \xi^{-Q}$, since the order of ξ does not divide Q . On the other hand we have $a_{k,i}^Q = 1$ for $i \neq l, l'$, since $a_{k,i} \in G_Q$ for $i \neq l, l'$. Thus Eve successfully identifies Alice's cryptotext if the cheating of M_1 and M_k is not detected.

Clearly, the values of b_1, c_1, d_1, e_1 , and f_1 are not changed by replacing L_0 with L'_0 in Step 1. Neither is a_1 , since $(\xi a_{0,1})(\xi^{-1} a_{0,2}) \prod_{i=3}^N a_{0,i} = \prod_{i=1}^N a_{0,i} = a_1$. Similarly, b_k, c_k, d_k, e_k , and f_k are not changed by replacing L_k with L'_k in Step 2. Neither is a_k , since $(\xi^{-1} a_{k,l})(\xi a_{k,l'}) \prod_{i \neq l, l'}^N a_{k,i} = \prod_{i=1}^N a_{k,i}$. Similarly as in the second attack of Section 3.3, we conclude that the proofs of knowledge produced by M_1 and M_k are deemed valid by the other mix-servers. If we assume that Alice and Bob are honest, their inner triples are never traced back and cheating is not detected.

If $\xi = \xi^{-1}$ Eve can only conclude that Alice sent a message from the set $\{m_l, m_{l'}\}$. This breaks the privacy of Alice, but Eve can also identify Alice's message uniquely by choosing Bob to be a corrupted sender.

Remark 5. Our attack may be viewed as a novel combination of the ideas in Lim and Lee [16] and Pfitzmann [24, 23] in that we use elements in $\mathbb{Z}_P^* \setminus G_Q$ to execute a “relation attack”. However, Lim and Lee [16] focus on *extracting the secret key*, and their attack requires that $P - 1$ contains a relatively large smooth factor. Their proposed countermeasure is to choose P such that $(P - 1)/2$ consists of large prime factors. This implies that the protocol leaks at most one bit. Our attack on the other hand use elements of $\mathbb{Z}_P^* \setminus G_Q$ more subtly, and is applicable⁴ and practical for *any* choice of P . In particular the type of P proposed by Lim and Lee.

Why the Attack Is Possible. The attack exploits that a mix-server M_j does not verify that all elements in its input L_{j-1} are in G_Q . M_1 uses this to “tag” elements in L_0 , which lets them be identified by the last mix-server M_k .

⁴ More generally, if the implementation of arithmetic for G_Q is induced from a larger group $G \supset G_Q$ such that $|G|$ is known and the protocol fails to verify that inputs are in G_Q and not only in G , then our attack is applicable.

3.5 Fifth Attack: One Corrupt Mix-Server

In Proposition 3 of Golle et al. [9] the authors claim that their protocol satisfies the following strong definition of public verifiability if there is a group G_Q in which the discrete logarithm problem is hard.

Definition 1. (*Public Verifiability (cf. [9])*) *A mix net is public verifiable if there exists a polynomially bounded verifier that takes as input the transcript of the mixing posted on the bulletin board, outputs “valid” if the set of valid outputs is a permuted decryption of all valid inputs, and otherwise outputs “invalid” with overwhelming probability. Note that to prove public verifiability, we consider an adversary that can control all mix-servers and all users.*

Unfortunately, Proposition 3 of [9] is false. This is seen as follows.

Suppose that there are 4 senders, and that the adversary corrupts two of the senders and the first mix-server M_1 . Let $(u_i, v_i, w_i) = (g^{r_i}, y^{r_i} m_i, H(u_i, v_i))$, and set $\alpha_i = ((a_{0,i}, b_{0,i}), (c_{0,i}, d_{0,i}), (e_{0,i}, f_{0,i})) = (E_y(u_i), E_y(v_i), E_y(w_i))$ for $i = 1, 2, 3, 4$. Then define $\alpha'_3 = ((a_{0,3}, b_{0,3}), (c_{0,3}, ad_{0,3}), (e_{0,3}, f_{0,3}))$ and $\alpha'_4 = ((a_{0,4}, b_{0,4}), (c_{0,4}, a^{-1}d_{0,4}), (e_{0,4}, f_{0,4}))$, for some $1 \neq a \in G_Q$. Suppose that α_1 and α_2 are sent to the mix-net by honest senders, and α'_3 and α'_4 are sent to the mix-net by the corrupted senders, with corresponding proofs of knowledge.

M_1 replaces α'_3 , and α'_4 with α_3 and α_4 . This does not change the value of the component-wise product $d_1 = v_1 v_2 v'_3 v'_4$ since $v'_3 v'_4 = v_3 v_4$, and the cheating is not detected, since α_3 and α_4 corresponds to valid inner triples, and thus not traced back. On the other hand the tuples α'_3 and α'_4 correspond to invalid inner triples (u_3, v_3, aw_3) and $(u_4, v_4, a^{-1}w_4)$. We conclude that the sets of valid inner triples in the input and output respectively differ, public verifiability is broken, and Proposition 3 of [9] is false.

Some may argue that this is not important since the adversary may only choose to correct invalid messages which she has previously prepared for this particular purpose. However, note that the adversary may *choose* whether to correct α'_3 and α'_4 . If she chooses not to correct invalid triples they are simply traced back and considered benign.

The following application shows the importance of this subtlety. We use the mix-net construction to run two independent elections (using different keys). First all votes for both elections are collected, and after some time both elections are closed. Then the mix-net is executed for the first election. Finally the mix-net is executed for the second election. In this scenario the adversary can insert specially prepared invalid triples (i.e. votes) in the second election, and then decide whether to correct these triples *based on the outcome* of the first election. This should clearly not be allowed, but may be acceptable in certain non-voting scenarios.

Why the Attack Is Possible. The attack exploits the fact that the first mix-server can choose whether to correct specially prepared invalid inner triples or not without detection.

4 Further Applications and Future Work

The attacks of sections 3.1, 3.2, and 3.3 all exploit the particular structure of the protocol of Golle et al. [9]. We have found no other protocol vulnerable to these attacks. In particular the protocol by Jakobsson [14] with similar structure is not vulnerable to the attack of Section 3.2.

The attack of Section 3.4 can be applied to break the privacy of Jakobsson [14], including the fixed protocol⁵ Mitomo and Kurosawa [18], and the attack of Section 3.5 breaks the robustness of Jakobsson and Juels [15]. These attacks are sketched in Appendices A and B.

An interesting open question is what other protocols are vulnerable to attacks based on malicious use of elements in $\mathbb{Z}_p^* \setminus G_Q$.

5 Conclusion

We have presented several practical attacks for the mix-net recently proposed by Golle, Zhong, Boneh, Jakobsson, and Juels [9], claimed secure by the authors. In particular we break privacy of any sender *without* corrupting any mix-server. Two of our attacks are easily adapted to break Jakobsson [14], including the Mitomo and Kurosawa [18] variant, and Jakobsson and Juels [15].

The attacks for mix-nets, presented here and in other papers [23, 23, 17, 4, 18], suggest that we should be careful with unproven claims of security for mix-nets.

Acknowledgements

I am grateful to Johan Håstad for his advise and support, and for providing a missing piece of the 4:th attack. My results would have been far fetched without the discussions I had with Gunnar Sjödin, Torsten Ekedahl, and Björn Grönvall. I also had discussions with Markus Jakobsson, Philippe Golle, Ari Juels, and Sheng Zhong. An anonymous referee pointed me to the work of Lim and Lee.

References

- [1] M. Abe, *Universally Verifiable mix-net with Verification Work Independent of the Number of Mix-centers*, Eurocrypt '98, pp. 437-447, LNCS 1403, 1998. 161
- [2] M. Abe, *Personal Communication*, april 2003, a paper will appear at ACISP '03. 161
- [3] D. Chaum, *Untraceable Electronic Mail, Return Addresses and Digital Pseudonyms*, Communications of the ACM - CACM '81, Vol. 24, No. 2, pp. 84-88, 1981. 160, 161
- [4] Y. Desmedt, K. Kurosawa, *How to break a practical MIX and design a new one*, Eurocrypt 2000, pp. 557-572, LNCS 1807, 2000. 160, 161, 165, 172
- [5] T. El Gamal, *A Public Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms*, IEEE Transactions on Information Theory, Vol. 31, No. 4, pp. 469-472, 1985. 161

⁵ We note that their proposal is given without security claims.

- [6] A. Fujioka, T. Okamoto and K. Ohta, *A practical secret voting scheme for large scale elections*, Auscrypt '92, LNCS 718, pp. 244-251, 1992. 161
- [7] J. Furukawa, K. Sako, *An efficient scheme for proving a shuffle*, Crypto 2001, LNCS 2139, pp. 368-387, 2001. 161
- [8] S. Goldwasser, S. Micali, *Probabilistic Encryption*, Journal of Computer and System Sciences (JCSS), Vol. 28, No. 2, pp. 270-299, 1984.
- [9] Golle, Zhong, Boneh, Jakobsson, Juels, *Optimistic Mixing for Exit-Polls*, Asiacypt 2002, LNCS, 2002. 160, 161, 162, 164, 166, 167, 171, 172
- [10] Golle, Zhong, Boneh, Jakobsson, Juels, *Private Communication*, 16 October 2002. 167
- [11] M. Hirt, K. Sako, *Efficient Receipt-Free Voting Based on Homomorphic Encryption*, Eurocrypt 2000, LNCS 1807, pp. 539-556, 2000.
- [12] M. Jakobsson, *A Practical Mix*, Eurocrypt '98, LNCS 1403, pp. 448-461, 1998. 161, 167
- [13] M. Jakobsson, D. M'Raihi, *Mix-based Electronic Payments*, 5:th Workshop on Selected Areas in Cryptography - SAC'98, LNCS 1556, pp. 157-173, 1998. 161
- [14] M. Jakobsson, *Flash Mixing*, 18:th ACM Symposium on Principles of Distributed Computing - PODC '98, pp. 83-89, 1998. 160, 161, 172, 174
- [15] M. Jakobsson, A. Juels, *An optimally robust hybrid mix network*, 20:th ACM Symposium on Principles of Distributed Computing - PODC '01, pp. 284-292, 2001. 160, 172, 174
- [16] C. H. Lim, P. J. Lee, *A Key Recovery Attack on Discrete Log-based Schemes Using a Prime Order Subgroup*, Crypto '97, LNCS 1294, pp. 249-263, 1997. 160, 165, 170
- [17] M. Michels, P. Horster, *Some remarks on a receipt-free and universally verifiable Mix-type voting scheme*, Asiacypt '96, pp. 125-132, LNCS 1163, 1996. 161, 172
- [18] M. Mitomo, K. Kurosawa, *Attack for Flash MIX*, Asiacypt 2000, pp. 192-204, LNCS 1976, 2000. 160, 161, 172, 174
- [19] A. Neff, *A verifiable secret shuffle and its application to E-Voting*, 8:th ACM Conference on Computer and Communications Security - CCS 2001, pp. 116-125, 2001. 161, 164, 167, 168
- [20] V. Niemi, A. Renvall, *How to prevent buying of votes in computer elections*, Asiacypt'94, LNCS 917, pp. 164-170, 1994. 161
- [21] W. Ogata, K. Kurosawa, K. Sako, K. Takatani, *Fault Tolerant Anonymous Channel*, Information and Communications Security - ICICS '97, pp. 440-444, LNCS 1334, 1997. 161
- [22] C. Park, K. Itoh, K. Kurosawa, *Efficient Anonymous Channel and All/Nothing Election Scheme*, Eurocrypt '93, LNCS 765, pp. 248-259, 1994. 161
- [23] B. Pfitzmann, *Breaking an Efficient Anonymous Channel*, Eurocrypt '94, LNCS 950, pp. 332-340, 1995. 160, 161, 164, 165, 167, 170, 172
- [24] B. Pfitzmann, A. Pfitzmann, *How to break the direct RSA-implementation of mixes*, Eurocrypt '89, LNCS 434, pp. 373-381, 1990. 160, 161, 164, 165, 167, 170
- [25] K. Sako, J. Killian, *Receipt-free Mix-Type Voting Scheme*, Eurocrypt '95, LNCS 921, pp. 393-403, 1995. 161
- [26] C. Schnorr, *Efficient signature generation by smart cards*, Journal of Cryptology, No. 4, pp. 161-174, 1991. 163
- [27] Y. Tsionis, M. Yung, *On the Security of El Gamal based Encryption*, International workshop on Public Key Cryptography, LNCS 1431, pp. 117-134, 1998.
- [28] D. Wagner, *A Generalized Birthday Problem*, Crypto 2002, LNCS 2442, pp. 288-304, 2002. 164

- [29] D. Wikström, *An Efficient Mix-Net*, Swedish Institute of Computer Science (SICS) Technical Report T2002:21, ISSN 1100-3154, SICS-T-2002/21-SE, 2002, <http://www.sics.se>, (An implementation was demonstrated during the Conference of the Swedish Research Institute for Information Technology (SITI), feb 7, 2002). 161, 166
- [30] D. Wikström, *How to Break, Fix, and Optimize “Optimistic Mix for Exit-Polls”*, Swedish Institute of Computer Science (SICS) Technical Report T2002:24, ISSN 1100-3154, ISRN SICS-T-2002/24-SE, 6 december 2002, <http://www.sics.se>. 161
- [31] D. Wikström, *Four Practical Attacks for “Optimistic Mixing for Exit-Polls”*, Swedish Institute of Computer Science (SICS) Technical Report T2003:04, ISSN 1100-3154, ISRN SICS-T-2003/04-SE, 25 february 2003, <http://www.sics.se>. 161

A Attack for “Flash Mix”

Jakobsson [14] presents an efficient mix-net. We assume familiarity with his protocol, and sketch how it is broken using the attack of Section 3.3.

The adversary corrupts M_1 and M_k . During the “second re-encryption”, M_1 “tags” two arbitrary El Gamal pairs in its input by multiplying their first components with ξ and ξ^{-1} respectively. Then the honest mix-servers perform their re-encryption and mixing. When the last mix-server M_k is about to re-encrypt and mix the output of the previous mix-server M_{k-1} , it identifies the “tagged” El Gamal pairs, removes the “tags” by multiplying the first components by ξ^{-1} and ξ respectively, and then finally re-encrypts and mix the resulting list honestly. After the verification of the “first re-encryption”, this breaks the privacy of some randomly chosen sender, if the cheating goes undetected.

Cheating is detected if one of two things happen; the adversary by chance chooses a “dummy element” that is later traced back through the mix-chain, or if M_1 or M_k fails to play its part in the computation of the “relative permutations” correctly. The first event is very unlikely since by construction there are very few “dummy elements”. Since the “tags” are removed by M_k , and both M_1 and M_k follow the protocol except for the use of the tags, it follows that the cheating is not detected. It is easy to see that the changes introduced by Mitomo and Kurosawa [18] do not counter the above attack.

B Attack for “An Optimally Robust Hybrid Mix Network”

Jakobsson and Juels [15] presents a hybrid mix-net. We assume familiarity with their protocol and sketch how it is broken using the attack of Section 3.5.

Suppose that there are four senders, and that the i :th sender forms a cryptotext $(c_0^{(i)}, \mu_0^{(i)}, y_0^{(i)})$ of a message m_i . The adversary corrupts the last two senders and modifies their cryptotexts as follows before they hand them to the mix-net. It replaces $y_0^{(3)}$ by $\bar{y}_0^{(3)} = ay_0^{(3)}$ by and $y_0^{(4)}$ by $\bar{y}_0^{(4)} = a^{-1}y_0^{(4)}$ for some $a \neq 1$.

The adversary corrupts M_1 . M_1 then replaces $\overline{y}_0^{(3)}$ by $y_0^{(3)}$ and $\overline{y}_0^{(4)}$ by $y_0^{(4)}$ and simulates an honest M_1 on the modified input. Similarly as in the original attack this does not change the component-wise product $P_0 = y_0^{(1)} y_0^{(2)} \overline{y}_0^{(3)} \overline{y}_0^{(4)} = y_0^{(1)} y_0^{(2)} y_0^{(3)} y_0^{(4)}$. The `VerifyComplaint` procedure is never invoked, since all cryptotexts are valid. Thus the cheating is not detected.

We conclude that the set of cleartext messages corresponding to the set of valid cryptotexts in the input differs from the set of cleartext messages in the output of the mix-net. This breaks the robustness, i.e. Definition 4(b).