

LLMix : Quantifying Mix Network Privacy Erosion with Generative Models

Vasilios Mavroudis¹ and Tariq Elahi²

¹ The Alan Turing Institute

² University of Edinburgh

v.mavroudis@turing.ac.uk, t.elahi@ed.ac.uk

Abstract. Modern mix networks improve over Tor and provide stronger privacy guarantees by robustly obfuscating metadata. As long as a message is routed through at least one honest mixnode, the privacy of the users involved is safeguarded. However, the complexity of the mixing mechanisms makes it difficult to estimate the cumulative privacy erosion occurring over time. This work uses a generative model trained on mixnet traffic to estimate the loss of privacy when users communicate persistently over a period of time. We train our large-language model from scratch on our specialized network traffic “language” and then use it to measure the sender-message unlinkability in various settings (e.g. mixing strategies, security parameters, observation window). Our findings reveal notable differences in privacy levels among mix strategies, even when they have similar mean latencies. In comparison, we demonstrate the limitations of traditional privacy metrics, such as entropy and log-likelihood, in fully capturing an adversary’s potential to synthesize information from multiple observations. Finally, we show that larger models exhibit greater sample efficiency and superior capabilities implying that further advancements in transformers will consequently enhance the accuracy of model-based privacy estimates.

1 Introduction

With the resurgence in popularity of Mix networks, due in part to the Snowden revelations, and upcoming real-world deployments like Nym [15], HOPR [23], and Elixir it is increasingly critical and timely to have assurances that mix design privacy expectations stand up to reality. Mixnets protect against powerful adversaries that have complete surveillance of all traffic links, thereby enhancing the security assurances beyond those provided by Tor. Indeed, many traffic analysis attacks that work on Tor, e.g. website fingerprinting, are not readily applicable in mixnets as the timing and communication patterns are robustly obfuscated by the *mixing strategies*. Furthermore, analytic solutions for larger mixnets are not tractable. Thus most designers and operators carry out empirical privacy measurements, including specific attacks that are critical to defend against, to determine the optimal mixing strategy and parameters for their use case [39, 15, 16, 21]. The two most common metrics used are: entropy and likelihood difference ϵ . They are both used to statistically summarize the leakage of the network or a single node. However, by definition they

do not combine observations from multiple mixing rounds and thus do not estimate the cumulative leakage over time (only the average), which is more realistic since a mixnet is expected to be used by a population over multiple rounds and a time horizon. Our work addresses this gap in the capabilities of existing tools, enabling operators to measure the resilience of their mixing strategies over several mixing rounds.

In particular, we introduce an automated design-agnostic approach for discovering *information leakage* in mixing strategies without the need for concrete attack realizations. This approach complements existing tools, and shifts the focus from the “prior” knowledge of the adversary and the specifics of the system details towards the privacy goals of the system and the objectives of the adversary. Central to this approach are: 1) our methodology for generating machine learning-compatible privacy-measuring tasks that model the goals of the adversary, and 2) the use of a transformer model as empirical privacy estimators (with regards to the aforementioned tasks).

We first introduce an encoding that represents and captures all relevant mix node transmission metadata as a stream of messages and represents them in a format that is compatible with modern large language models. Using [24] we translate the high-level privacy goal, specifically sender-message unlinkability, into a privacy game and model it as two machine learning (ML) tasks (i.e. distinguishing between two senders, and one sender amongst many). We then generate network traces for a range of mixnet configurations and train the transformer model on them. During training the model learns the underlying rules of mixing by trying to generate valid mixnet traffic traces (i.e., guess the next message transmission). This traffic-aware model is then used to solve our privacy-estimating tasks given a new mixnet traffic trace.

Our model LLMix³ is a transformer (two variants) trained *from scratch* to process and classify mixnet traffic. We evaluate various mixnet strategies and find that configurations that impose the same average latency are not always equally robust with regards to the unlinkability property (Section 7.2). We focus on the leakage of individual nodes as our goal is to provide a best case (for the defender) privacy comparison of the different mixing strategies and parameters, without the added complexity of network topology choices. This setup is in keeping with the standard *anytrust* assumption where even one mix node must be capable of ensuring the privacy of the users [50, 26]. We also train a larger variant of our model to study if it benefits from an increased observation window. We confirm our hypothesis and show there is a clear relationship between the number of messages captured and privacy loss (Section 7.4). Comparing these results with classic statistical tools [39, 15, 21] we show that our proposed technique provides a better privacy estimation (i.e. the privacy information leakage). Finally, provide initial evidence that the number of *learnable* parameters directly influences the model’s sample efficiency (Section 7.3).

Overall, this work realises a task-to-model framework enabling ML advances to be applied to harden mixing networks and guide parameterization. In particular, we:

- Introduce a new traffic analysis task format that is solvable by language models and train from *scratch* a mixnet traffic analysis transformer model.

³ We will open-source our trained generative models, tools and scripts upon publication.

- Provide, to the best of our knowledge, the first design-agnostic generative model (LLMix) that accounts for cross-round leakage in mixing privacy estimations.
- Evaluate a range of configurations and find that some combinations achieve better privacy with the same latency.
- Study sample-efficiency, and show that larger models are likely to have better privacy-estimating capabilities.

2 Background

2.1 Mix Networks

First proposed as a mechanism for untraceable electronic mail [7], Mix networks (or mixnets) have since been widely adopted for several applications including secure e-voting [20], anonymous routing [9] and anonymous messaging [39, 1]. HOPR [23] and Nym [15] are two recent examples of real-world deployments using the latest advances in mixnet techniques as stronger alternatives to more established solutions like Tor and VPNs.

To transmit a message over a Mixnet, the sender selects a route over a number of nodes (or hops) before the message arrives at the recipient. A single honest mix node in the route ensures the unlinkability of the message. Each message is padded to a constant length before it is sent and then cryptographically processed at each hop (either decrypted or re-encrypted depending on the scheme used) to prevent traceable bit patterns, with an additional delay at each hop. The particular *mixing strategy* dictates this delay and is what differentiates mixnets from Tor (uses FIFO routing). At the expense of additional latency, the delay ensures that multiple messages from (ideally) different users are co-resident in each node thus creating confusion for the adversary to trace the path of messages through the network, thus preventing the linking of sender and recipient. Dummy messages can further obfuscate the link between sender and recipient at the expense of bandwidth.

Mixing nodes are the building blocks of mixing networks as each of them individually ensures the realisation of the mixing strategy. In fact, a mix net with a single node would provide the best possible mixing of the traffic available. However, for scalability and to avoid single points of failure, most modern mixnet designs use a stratified topology where nodes are arranged in ordered layers with messages traversing the network across the layers in sequence. This fragments the traffic thus reducing the homogeneity of the mixing. Most of the mixnets using such a topology operate under the anytrust assumption[50] where at least one server in the user’s path must be honest. Thus security comes from distributing trust across many relay operators.

2.2 Transformer Models

Many different approaches to language modelling have been proposed [45, 28, 35, 36], however massive scaling (i.e., large) in language model parameters has recently yielded unprecedented performance improvements across several tasks [14, 6, 41, 51, 52]. Large language models (LLMs) typically use a *transformer* neural network

architecture [47] designed to process sequential data such as natural language based on self-attention. Self-attention is an attention mechanism [3] that allows dependencies in sequential data to be modelled independently of their distance in the input and output. Importantly transformer models are highly parallelisable, enabling the scale necessary for LLMs. In what is termed “autoregressive training”, LLMs for natural language tasks are usually trained to predict the next word in a given sentence. This process is performed iteratively, generating one word prediction in each step. Despite the relative simplicity of this process, autoregressive training is sufficient to capture much of the syntax and semantics of language. LLMs can generate coherent and contextually relevant natural language, allowing them to perform well in many previously unsolved tasks [13, 30]. They are given a vocabulary set that defines all of the unique tokens the model can recognise and generate. Depending on how the vocabulary is defined, each token may correspond to a whole word, a subword, a character, or a byte. Tokenization splits raw text (e.g., a phrase, sentence, paragraph, or a document) into individual tokens from the vocabulary, used as input to the LLM. While transformers are generative models, they can nonetheless be used for classification tasks while also taking advantage of their ability to resolve long-range dependencies in sequential data. Section 4 outlines how we train them from scratch to process non-human language i.e., mixnet network traces.

3 Threat Model

We assume the standard global passive adversary (GPA) that is able to observe all network traffic between users and mixnode under examination [39, 7, 2, 21]. The adversary observes a fixed number of network events i.e. messages entering and leaving the network. The messages are all indistinguishable from each other (Sphinx [12]). The adversary does not actively inject, drop, replay or delay messages, and does not operate any compromised end-users, i.e. sender or recipient of a message. The adversary has the ability to corrupt nodes, however, at least one of the mixnodes each message is traversing through is assumed to be honest. This is the anytrust assumption and is commonly used by many known designs (e.g. Vuvuzela [46], Karaoke [25], Loopix [39], Nym [15]). Corrupt nodes are assumed to operate on a FIFO manner as the adversary can fully observe their operation and hence they do not contribute to the mixing [39].

3.1 Adversarial Goals & Tasks

The overall goal of an adversary is to breach the privacy of the users, specifically to break the *unlinkability* privacy property of the mixnode and link the communication between the sender and recipient of the message. We focus on this goal since it is fundamental to mix networks and the basis of other privacy goals. Kuhn et al. formalize this *sender-message unlinkability* privacy property as the $SM\bar{L}$ notion, a fundamental desired property shared by many mixnet designs and thus of main concern when evaluating mixnet designs and configurations (e.g., Loopix [39]). We choose to focus on this property due to its importance in most mixnet use cases.

However, our methodology can be easily adapted to support any notions defined in [24].

Our work extends the foundational work of Kuhn et al. by translating their established privacy goal and concept framework into corresponding tasks within the realm of machine learning. We then use the success of the adversary at the ML task to gauge the privacy level of a given mixnet with respect to that goal. Mixnet designers and operators can use this information to then make design and deployment decisions such as rule out mixing strategies and configurations that provide subpar privacy. Note that in the literature the privacy “game” commonly involves guessing the sender of a *single* message rather than identifying a *persistent contact*. We consider the latter to be more pragmatic. In the former case, when the privacy leakage is minor, the measurement noise will not indicate that the adversary has gained any substantial advantage by observing a single message. However, leakage is accumulative and observing several messages may result in an advantage eventually, as our results show. We now introduce the task that instantiates $SM\bar{L}$:

Task: One of Two

Given a recipient B , the adversary aims to identify the user A who sends messages to B . The network has several actively communicating users but the adversary has to choose between *two* potential senders (e.g., due to prior knowledge through external information). A coin-flipping adversary has a $\frac{1}{2}$ chance of guessing correctly in this task and the privacy loss is given by any adversarial performance that exceeds this mark.

4 Converting Privacy Goals into ML Tasks

As discussed in Section 2, transformers were designed to solve natural language processing and computer vision tasks. We argue that network traffic exhibits properties and structure that is analogous to those applications, thus making network traffic a suitable application area. We now delve into the details of leveraging the capabilities of an LLM to estimate the privacy of a given mixnet with regards to a specific adversary (i.e., desired privacy property and corresponding ML task).

4.1 Network Traffic Encoding

Modern mixing strategies and the use of cryptographic message formats (e.g., Sphinx [12]) leave only a few types of metadata exposed. Specifically, for each message transmission event, an observer knows either the sending or receiving systems/parties, the direction of the transmission, and the relative order of this transmission to the rest of the transmissions taking place in the network. Figure 1 shows how an adversary can represent a snapshot of the network’s activity as a sequence that encodes the aforementioned metadata.

In each case, the sequence produced retains all the information that a passive adversary can collect by eavesdropping the network links and is in a format compatible

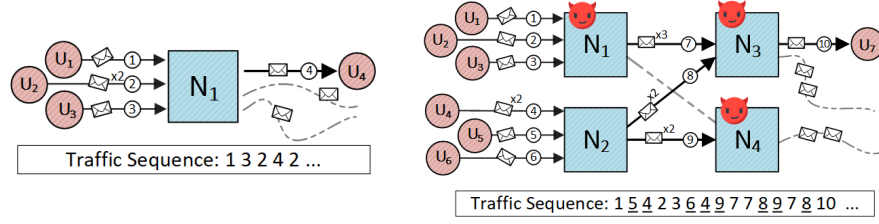


Fig. 1: Intercepted mixnet traffic can be efficiently represented as a sequence. Messages are replaced by the numerical id of the network link they traversed. The leftmost figure shows the traffic as it is routed through a single node, while the rightmost one shows a mixnet with several compromised nodes (anytrust assumption). The underlined link ids correspond to traffic that they adversary does not have full visibility on due to N_2 being honest.

with LLM input sequences. This is because we have translated network traffic into a language with the following characteristics: Each network event is represented by its corresponding link id. Each “word” (i.e., link id) is unique and has a single meaning. Two or more words can be combined to form a sentence (i.e., a traffic sequence) where the relative order of the events matters.

Given these characteristics and the ability of transformers to process natural language, we argue that the network activity “language” outlined above is considerably easier for an LLM to learn (compared to a natural language). Our encoding ensures there is no lexical ambiguity (i.e., one word with two meanings) and the vocabulary is relatively small for networks of moderate size. Transformers are mathematical functions and thus the input sequences must be converted into numbers. For natural languages, tokenizers are used to convert character-based words to integers. In our case, the mapping is even more straightforward as our link ids are already monotonically increasing integers (a word-level tokenizer). A “0” represents the absence of network activity.

Note that in NLP the encoding used to represent the information is always the same (e.g., voice or written text of a human language) regardless of the task at hand (e.g., classification, sentiment analysis, classification). This is an important characteristic of the way tasks are modelled to be solved by LLMs and part of why LLMs generalize so well. In the same vein, our network activity “language” is also generic as it represents the activity taking place in the mixing network in a lossless manner regardless of the end-goal of the adversary. This approach is unlike classic ML methodologies that derive features specific to the objectives of the analysis that the model then operates on (e.g., website fingerprinting [48]).

4.2 Privacy Properties & Games

In general, we fix the set of security or privacy properties a mix node will be evaluated against. As mentioned before, Kuhn et al. [24] provide a comprehensive framework that formalizes an anonymous network’s privacy goals as *notions* and

defines a hierarchy between them. Given a particular use-case, the framework allows practitioners to define the communication setting and express their intuitive privacy goals as formal privacy notions and *games*. Using this methodology, the steps of the game that corresponds our task above (Section 3.1) are: 1) the adversary picks two potential senders and a recipient from the mixnet’s population; 2) the challenger checks if the senders and the recipients are distinct; 3) based on a random bit b the challenger inputs a scenario into the mixnet testbed; 4) the adversary observes the traffic sequence and outputs a ‘guess’ b' as to which scenario was executed.

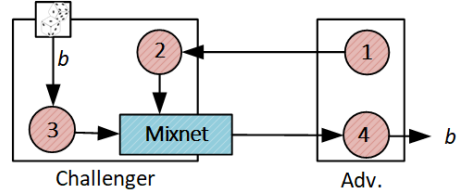


Fig. 2: A sender indistinguishability game: 1) the adversary picks two potential senders and a recipient; 2) the challenger checks the selections; 3) the challenger randomly chooses and runs the scenario; 4) the adversary observes the traffic sequence and outputs a ‘guess’ b' .

Once the game(s) has been defined, we use a simulator, or an emulator, of the network to generate data from the scenario. From each run, we get the traffic sequence and ground-truth label (i.e., bit b above). The label is 0/1 depending on if user A or user B was the persistent contact. We can then compile a dataset (each row has a traffic sequence and its label) and train models.

5 Data Collection

Before we present our experimental results, we discuss our setup and the rationale for our design decisions for our testbed. The testbed realises mixing nodes and implements the threshold, pool, and Poisson mixing strategies. Our testbed is implemented in Python 3.12 and is comparable with those used in the mixnet literature (e.g. [21, 37, 39]).

5.1 Parameters

The number of active users (N) is at least 3 (i.e., two senders, one recipient) but in practice, we study significantly greater user numbers. The rate of sending r_i per user is a ratio i.e., if $r_i = 0.5$ then the user sends 0.5 messages in every unit of time. The global sending rate $R \leq 1$ (up to 1 message every time unit). This means that time moves forward only when there is a message transmitted somewhere in the network. In other words, the adversary will simply discard periods of inactivity from

their dataset. In cases where the users all behave uniformly, $R = r_i * N$. The average end-to-end latency is denoted l and is not capped as it depends entirely on the mixing parameters λ , n and p depending on the mixing strategy used. Note that p is a ratio between the messages held in the pool and n . p is always smaller than 1 as the pool cannot be of equal (or larger) size to the threshold buffer. Finally, the users need to have at least one contact otherwise they are not considered active and should not be accounted for in N .

5.2 User Traffic

Mixing networks rely on user traffic to provide anonymity. Consequently, a large number of users and hence more traffic makes it easier for defenders to be assured of a high level of privacy [18].

Threshold mixes have been typically studied on the assumption that all the users participate equally in each *batch* of messages shuffled together. As stated in the original paper introducing threshold mixes [8]: “...each participant supplies the same number of messages to each batch”. Having perfectly balanced *batches* where all the users are equally represented, significantly strengthens the privacy of the network but the usability impact of such a rule would make threshold mixes completely impractical to use. Instead, and more realistically, we assume that all active users send messages at approximately the same rate. This is favourable to the defender as it maximizes the sender diversity for the messages traversing the network and prevents users from trivially standing out. It also allows us to evaluate the quality of the mixing while minimizing the noise that a skewed user traffic distribution would introduce. In practice, real-life sending patterns are not perfectly uniform and thus the adversarial performance reported in Section 7 should be treated as an upper-bound for the privacy a given configuration can guarantee. Moreover, sophisticated (active) adversarial strategies can exacerbate uneven sending distributions further so as to erode privacy even more.

We consider only two-party communication (cf. multicasting) via *messages* at the application-layer. We believe that this is representative of two-party exchanges that take place in short bouts of messaging activity. The *contact* of each user is chosen by sampling uniformly from the user list and does not assume any reciprocity (i.e., if A elects B as its contact, then B does not necessarily elect A too). This is to avoid introducing priors into the communication graph.

Focusing on high-level messages, rather than network-layer packets, abstracts away from heavy hitters who send larger texts, pictures, videos (that require multiple TCP/UDP packets) from standing out. In practice, cryptographic formats such as Sphinx [12] and techniques such as dummy packets prevent trivial deanonymization of heavy hitters, however some risk remains.

Consistent with our threat model we do not consider active attacks such as the $n - 1$ message attack [43]. Such attacks provide a significant advantage to the adversary but are more detectable and thus riskier, compared to our passive adversary.

5.3 Burn-in Phase

Before each testbed run, we perform a burn-in that aims to prefill the buffers and the pools of the mixing nodes. This ensures that the adversary gets to observe traces from a mixnet that is properly initialized. Poisson mixes are particularly sensitive to poorly initialized networks. To find the minimum time it takes for the buffers to stabilise, we tracked the number of messages in the mixnode buffers over time and for different λ values. We found that after 300 timesteps all buffers are stably populated for $\lambda \leq 50$. We further corroborated the stability of the buffer and the mixing quality by measuring the buffer’s entropy for over 1,000 timesteps. Single-node threshold mixes without a pool do not require burn-in but we followed the same approach for threshold mixes with a pool. Based on our findings, we set each burn-in to last 4096 timesteps followed by another period of random duration (between [1,4096]). This random component ensures that there is no alignment (e.g., traces always start after 96 messages have been transmitted from the network’s exit node) between the traces generated. During this period, the “suspect” senders do not send any messages.

6 Experimental Methodology

For our experiments we use the Longformer [4] and train from scratch two variants (with moderate and large parameter number respectively). This is an important distinction as all pretrained large language models (LLMs) found online are trained on natural language tasks and do not transfer to network traffic.

Moreover, some of these model architectures (e.g., LLaMa) are very capable but also require an extremely high cost investment in compute to train them. We opted for the Longformer as it introduces a variation of self-attention (i.e., local attention) that scales linearly to the sequence length. It has approximately 149 million parameters and its primary advantage is that it supports considerably longer input sequences at a relatively low computational cost. In contrast, other self-attention mechanisms grow quadratically with the sequence length, making it infeasible to process long sequences. With local attention every token attends only to tokens in its vicinity defined by a window w . In particular, each token attends to the $\frac{1}{2}w$ preceding and the $\frac{1}{2}w$ succeeding tokens. Longformer combines this with *dilation* allowing for increases in the window size without incurring additional memory costs.

Despite its sophistication, the Longformer is a generative model and simply predicts the subsequent token based on the past. To use an LLM for classification, we follow the standard practice and replace the first token of every sequence with a special classification token “[CLS]”. The final hidden state corresponding to this token is used as the aggregate sequence representation. This representation is then passed through a simple classifier (usually a Linear layer) to get the predicted class label. This approach was first introduced in BERT [14] and is an established way of training an LLM for a classification task.

7 Experiments

In this Section, we use our model LLMix to evaluate the mixing strategies and parameters under various circumstances and configurations. In all our experiments, we use our testbed to generate data from a network with a specified number of users, topology and mixing parameters. We conduct our experiments under the anytrust assumption which requires that as long as a message crosses one honest node its privacy should be preserved (see also Section 3). We split the produced labelled data into three datasets: *training*, *validation* and *testing*. All our reported results correspond to the models’ performance on the testing dataset. Testing dataset samples were not encountered during training or when tuning our models’ hyperparameters (validation dataset).

7.1 Configuration & Hyperparameters

LLMix is a variant of the Longformer model with its parameters optimized for traffic analysis tasks and lower computational costs. We now describe the most important parameter choices and our rationale. Note that many of the considerations and the practical rules of thumb are specific to LLMs.

Hardware. We use 1) 4 DGX-MAX-Q Nodes each with 8x Nvidia V100 GPUs with 32GB of VRAM, and 2) an HGX100 GPU planar with 4x Nvidia A100 80GB GPUs.

Software. We use Python 3.10.4 with the technically relevant packages being: PyTorch 2.0 backend, Transformers 4.30.2 [49], Accelerate 0.23.0⁴ and NumPy 1.26.0.

Batch size. The batch size governs the training speed of the model. Very small values (e.g., 1,2) can result in extremely slow training that will not converge within a reasonable timeframe. Thus larger values (e.g., 256, 512, 1024 or larger) are preferable as they increase the rate of training and prevent instabilities during training [27]. Note that the relationship between the batch size and the training speed is not linear. The maximum batch size supported by the V100 GPUs (i.e., GPU memory) was 32. This value, although moderate, is consistent with other works training Longformer models [4].

Gradient Accumulation. Gradient accumulation splits the batch of samples into mini-batches that are processed sequentially. While the mini-batches are processed, the gradients of those steps are collected and accumulated without updating the model. The model is updated once all the mini-batches have returned. We use 8 gradient accumulation steps that bring the effective batch size to 256 (32x8). Understandably, this process introduces significant latency in the training process as it requires moving data in and out of the GPUs VRAM. However, it significantly improved the training stability of our models and minimized the instances of non-learning models when increasing the complexity of the task in curriculum learning.

Epochs. We did not define a set number of epochs but instead let the models train at each level of the curriculum until there was no significant evaluation accuracy improvement (1% or more) for at least 20 epochs. This is standard practice and is

⁴ <https://huggingface.co/docs/accelerate/index>

implemented using an *EarlyStoppingCallback* that takes as parameters the number of epochs (20) and the minimum improvement difference (0.01).

FP16 Training. We used half-precision floating-point numbers (i.e., FP16) as it allows for faster computations. The main speedup comes from storing the layer activations in 16-bit precision and thus making numerical operations easier. This also results in 16-bit gradients. However, this does not provide any memory usage benefits. This is because the model is saved in the GPU memory in *both* 16-bit and 32-bit precision and can thus even use more memory than a non-mixed precision version. Moreover, the gradients are converted back to full precision for the optimization step thus preventing any potential memory savings. We experimented with various batch sizes and precision settings and in our setup FP16 allowed for faster computations without impacting the maximum batch size we could fit in the GPU’s VRAM.

Attention Window. We set the attention window (w) to 128. The motivation for this was that we wanted each attention head to be able to observe a full message shuffling round if possible. In this case, the window allows for up to 64 tokens observed before and 64 after the mixing. Thus a mixing round on a mix with threshold 50 will be fully observable. Of course, LLMs employ multiple attention heads and combine observations from two or more heads if needed. We validated this in our use case by setting the window to 64. The model was still able to analyse the traffic even when the threshold was larger than $\frac{1}{2}w = 32$. Moreover, preliminary experiments with threshold values that exceed the 128 token window also confirmed this.

Vocabulary Size. This defines the maximum number of different tokens that can be represented as tokens in the LLMs’ input. We used a vocabulary size of 10,000 tokens as this was enough for the configurations we evaluated. In particular, the vocabulary size needed is about twice the number of network links.

Model Architecture. Due to hardware constraints we also had to reduce the dimensionality of the encoding and the pooler layers to 128 (i.e., hidden size). A common value is 512 when processing longer natural language sequences. As well as, the number of hidden layers to 8. Values between 8-12 are common for smaller models while values of up to 30 are used for larger ones. Finally, we set the number of attention heads to 12 which is a commonly used value. We arrived at this combination of values through a grid search of the parameter space, trying to balance learning performance and speed within the capabilities of our hardware.

7.2 Latency vs. Privacy

We now focus on the three most common mixing strategies from the literature: Poisson mixing, Threshold mixing and Threshold mixing with a pool. These methods all try to achieve the same goal (privacy) but their usability impact (latency) is not identical. In this experiment, we study the level of privacy achieved by a single node implementing one of the above mixing strategies with regards to the average latency imposed to the messages traversing the network. Using a single node, allows us to isolate any impact network’s topology might have to the results. The single node serves could also be considered as a high level abstraction of a whole mixnet where the adversary does not have access to intra-node communication links.

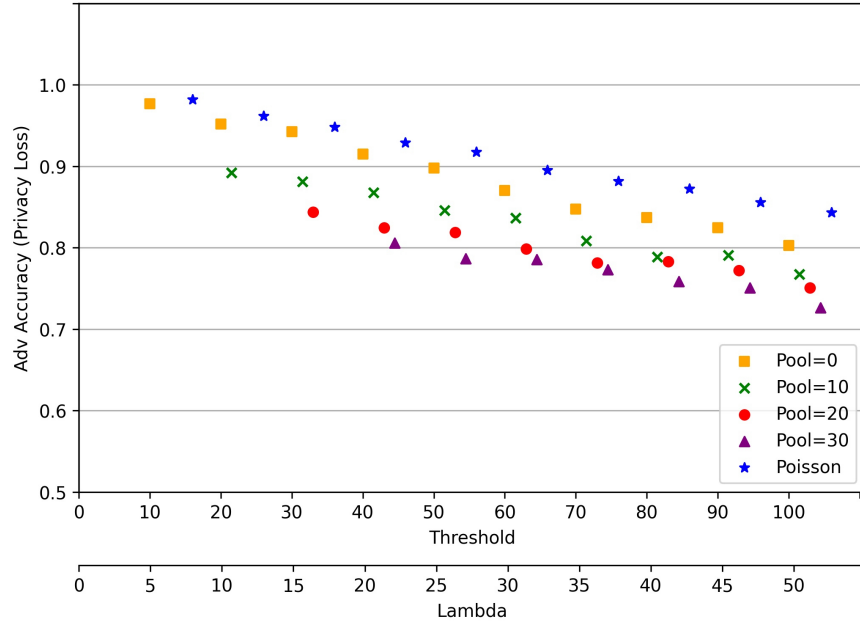


Fig. 3: Privacy-loss estimation for different types of mixing strategies. The estimation task is guessing the sender of a message arriving to a given recipient. The total users are 100 and the possible senders are 2. The x axes are aligned by the latency each threshold/lambda value imposes. Lower values are better as values closer to 0.5 indicate better privacy. A node with a pool=0 is equivalent to a threshold-only node.

We evaluate mixnodes with thresholds in the range of 10-100 and the corresponding (with regards to the latency imposed) λ for the Poisson nodes in the range 5-50 (with a step of 10 and 5 respectively). Our network has 100 users who each generate 1 message per 100 seconds. The privacy-estimation task assumes an adversary who given a recipient r guesses r 's persistent contact between two possible senders (Section 3.1).

We started by training one LLMix model on labelled data collected (2048 messages) from a mixnode with threshold 10 and no pool. This initial training occupied a V100 GPU and took approximately 2 days before the LLM performed measurably better than a random guesser. We kept training the network for another 2 days until it converged to its final performance as reported in Figure 3. Then we took advantage of the homogeneity of the datasets and used this trained model as a pre-trained basis for follow up training (curriculum learning). We trained another 43 models for all the remaining mixing strategies and configurations (Figure 3).

We observe that Poisson mixing performs worse than every other strategy and configuration. This can be explained by the way *continuous* mixing works. In particular, batch based mixing strategies (e.g., threshold) do not flush messages out unless there are enough messages in the pool. In contrast, continuous mixes favour usability (e.g., to prevent messages from staying impractically long in the node) and treat every message independently. On average these two approaches ensure that each messages is shuffled with at $n - 1$ other messages when in the node. However, there are edge cases where the mixing quality deviates from the average. As explained in Section 3, we assume a uniform sending behaviour for all the users which is a favourable assumptions for the defender.

The number of messages in the buffer exhibits substantial variance which is the primary cause of the model's lower privacy estimate. Moreover, we also observe that a pool provides a significant advantage to a threshold mix. To better evaluate the cost of a pool in the usability of the network, we measured the latency of all the configurations in evaluation. Figure 4 shows that a mix node with a threshold $n = 70$ and a pool holding 10 messages $p = 14\%$ (setup A) provides about the same privacy with a threshold $n = 100$ and no pool (setup B). Interestingly, setup A also imposes higher average latency to the messages traversing the network while the increased latency variance introduced by setup A is less than that of B. Furthermore, a node with $n = 80$ and a pool holding 10 messages $p = 13\%$ offers more privacy than setup A, lower average latency and comparable latency variance. Note that the latency in Poisson nodes equals the λ value, whereas the latency of threshold mixes is equal to approximately $n/2$. This is because messages that arrive in the buffer later have to wait less time before the threshold is flushed. In fact the first message will wait n seconds (assuming an arrival rate of 1 second/message) while the n 'th message will not experience any delay.

Overall, the trade-off curves (privacy vs. latency) for each mixing configuration differ and thus a privacy evaluation is useful to identify the setups that strike the best balance. Moreover, none of these configurations provides perfect privacy under the anytrust assumption (see Section 3). Note that, if we assume a weaker adversary that does not have access to intra-node network links, a mixing network can be abstracted

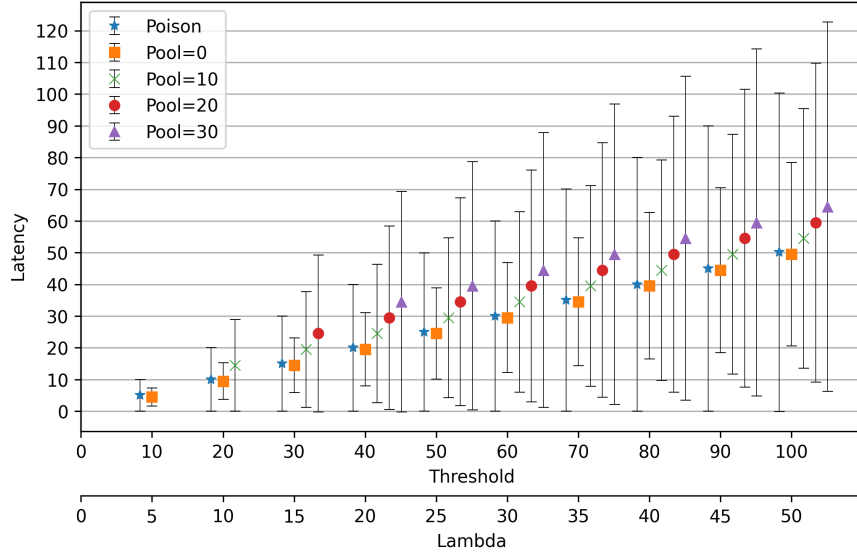


Fig. 4: Mean latency and standard deviation (in seconds) for various configurations in Threshold, Pool and Poisson Mixes. We observe that different configurations might have the same mean latency (variance might differ). A node with a pool=0 is equivalent to a threshold-only node.

as a single mixing node as long as the correspondence between the single node and the multi-node setup mixing is accurate.

7.3 Num of Observations & Model Size

One of our hypotheses is that advancements in the capabilities of transformer models will translate directly into better privacy estimations. To test this hypothesis we now train and evaluate a larger model. We use the same setup with Experiment 1 (Section 7.2). To train the larger model we used Nvidia A100 80GB GPUs as the increased number of parameters did not fit in the original V100 with 32GBs of VRAM. However, as the models grow small batch sizes do not suffice for the model to learn. We thus increased the batch size to 200 but disabled the gradient accumulation as we found that it slowed training substantially. The increased attention window allows one head to observe more network events which together with the larger number of observations (4096 vs 2048) could give to the model greater inference capacity. Note that by increasing the observation length of each sample, we double also the space it occupies in the GPU's memory.

Table 1 shows the accuracy of this model for sequences with various observation lengths. However, trained models require a fixed input length and will not work if samples of a different size are provided. As our larger model was trained on samples of 4096 events, we have to provide compatible samples. To achieve this, we used the

same testing dataset (never seen during training). For each sample in this dataset, we randomly selected a region with the size we wanted to test against, and masked the rest of the events with zeros ('0' represents no network activity).

From Table 1, we observe that indeed the larger model is more capable than the moderate sized-one. In particular, for a sequence of length 2048 the larger model achieves an accuracy of 0.88 while the smaller model reaches an accuracy of about 0.813. This shows that not all the performance gain is due to the increased observation length. However, longer sequences indeed allow the model to do even better (0.95). In fact, the accuracy of the model depends heavily on the number of messages sent by the targeted sender (or observation length), and the performance declines as the sequences half in size. Interestingly, even 1 message sent allows the model to guess correctly with odds considerably better than random (accuracy 0.58).

Sequences		Metrics		
# Obs	# Messages	Adv. Advantage (ours)	Likelihood Diff.	Entropy
4096	20.7	0.958*	0.269	5.824*
2048	10.9	0.884*	0.270	5.833*
1024	5.1	0.776*	0.270	5.849*
512	2.0	0.661*	0.260	5.865*
256	1.0	0.583*	0.262	5.859

Table 1: The relationship between the observations of the adversary, the messages sent by the real sender to the target recipient, the estimates of our metric, as well as the entropy and the likelihood difference ϵ measured. The network has 100 active users and a threshold of 100. Adversarial advantage values closer to 0.5 indicates better privacy. We use * to mark the mean values that are statistically significantly (p-value 0.05) different from their adjacent means (rows above and below).

7.4 Comparison with other metrics

Entropy and likelihood difference ϵ have emerged as standard metrics to quantify the adversary's advantage [42, 17, 21, 39, 38]. We now show that these measures have limitations and discuss how our proposed model provides more accurate estimations.

The entropy metric quantifies the information contained in the anonymity probability distribution of possible senders and receivers, of a given message, as viewed by an adversary. The effective anonymity set size is defined as the entropy of the anonymity probability distribution, which can be interpreted as the number of additional bits needed to uniquely identify the specific sender or receiver of a chosen message. Compared with the anonymity set size, entropy provides a probabilistic measure of the information that an adversary can determine from a sequence. More formally, let X be a discrete random variable over the finite set X with probability

mass function $p(x)=\Pr(X=x)$. The Shannon entropy $H(X)$ of a discrete random variable X is defined as:

$$H(X) = - \sum_{x \in X} p(x) \log p(x) \quad (1)$$

The likelihood difference ϵ [39] provides the expected difference in likelihood that a message leaving a specific terminal mix node is sent from one sender in comparison to another. In comparison to entropy, likelihood difference ϵ focuses on the ratio between the probabilities that a selected message was sent by one of two chosen senders. It thus allows anonymity to be quantified more pessimistically for a strong adversary who is assumed to know a priori that one of two specific senders are communicating with a certain recipient. Given a message, its likelihood difference is calculated as: $\epsilon = |\log(p_0/p_1)|$, where $p_0 = \Pr[U_0]$ and $p_1 = \Pr[U_1]$ denote the probabilities of users U_0 and U_1 being the senders of that message respectively.

Continuing with Table 1, we examine the likelihood difference between sequences with a varying number of messages from the real sender to the target recipient. The mean values shown in the table correspond to 5 rounds (with 256, 512, 1024, 2048, 4096 observations) of 1,000 sequences each. Each metric should capture the additional privacy loss as the number of messages increases. However as seen in the table, the difference between the 256-512, 512-1024, 1024-2048, 2048-4096 is not statistically significant and thus the likelihood difference ϵ fails to capture the extend of the leakage. Moreover, the difference between 256-2048 was not statistically significant either, while the reported values are very close to what [39] considers safe. In comparison, our model shows that there is a clear leakage in each of these cases and all mean differences were statistically significant.

The entropy metric captures some of the leakage but does not reveal its full extend. To understand why, we use an example from [21]. An entropy of 10 bits for a mixed message arriving at a recipient, indicates that the message is “as anonymous as if it was perfectly indistinguishable among about a thousand ($2^{10} = 1024$) other messages”. If the same user sends a follow up message that mixes in the exact same way with messages from the same recipients, there is no additional privacy loss. However, in practice the subsequent message will be shuffled with a different mixture of messages from different senders. Even if the entropy of the subsequent message is also 10 bits, the probability distributions of the two messages will differ e.g., allowing the adversary rule out combinations of senders and recipients (something our model is capable of). Subsequent messages will provide additional information that will allow further inferences. In our case all the entropy measurements were between 5.824 (for 4096) 5.859 (for 256) which corresponds to the messages being indistinguishable between 56.64 and 58.04 other messages respectively (at most $1/56.64$ changes of guessing correctly the sender of a message). Entropy thus captures some of the leakage but as our model shows it under-reports its full extent. Interestingly, all the mean differences (except for 256-512) were statistically significant. The reason for this that the are only 1 and 2 leaky messages respectively and consequently only 1-2 low entropy measurements that did not provide enough signal.

Overall, we conclude that while statistical tools are easy to use for debugging and quickly iterating when parameterizing a node, they fail to capture the full extent of

the leakage. This gap in capabilities can be addressed by our proposed methodology and model.

8 Related Work

The problem of traffic analysis for privacy purposes remains open and is a relatively popular theme in anonymous communication networks research. However, there are only a few sub-areas where ML has been actively used so far.

Website fingerprinting attacks against the Tor have employed standard machine learning techniques for classification such as k-NN [48], Support Vector Machines [34], random forests [22], and more recently convolutional deep neural networks [44, 5]. In webpage fingerprinting, Danezis et al. [11] used a Hidden Markov Model to fingerprint pageloads from a static dataset. Miller et al. [29] fit a Gaussian distribution on their dataset while Dubin et al. [19] used deep learning to fingerprint traces from video streaming services.

Attack-agnostic traffic analysis for mixnets using modern ML models remains an underexplored area as tools and techniques from other types of anonymity networks are not applicable. Past works have focused on studying specific attacks against specific mixing strategies [32, 40, 33], covering the range of different mix types from threshold [40], to pool [40], to continuous [10], and consider enhancements such as dummy traffic as well [31]. Other works aim to provide optimal parameters to secure continuous mixnets [21, 37] like Loopix. This work can be viewed as an extension to these tools and can be used in conjunction with them. This has the potential to more fully cover relevant sources of information leakage that the more limited measurements in the related work might miss.

9 Conclusion

We introduced a framework that allows system designers and operators to estimate the privacy of their mixing strategies against practical adversaries. Our approach defines a “language” and a task format that are compatible with modern large language models and thus makes it straightforward for all advances in language models to be directly incorporated in privacy estimation. We believe that such estimations are close to the actual advantage of a sophisticated adversary, and as models advance, we expect the gap to become even narrower. While the list of adversarial tasks can not be exhaustive, privacy goal systematization efforts (e.g. [24]) can guide operators towards identifying and “taskifying” their primary privacy objectives. With this framework in place, it is now straightforward for designers and operators to evaluate their configurations against the latest language models, under conditions and assumptions that fit their use cases.

References

1. Nikolaos Alexopoulos, Aggelos Kiayias, Riivo Talviste, and Thomas Zacharias. MCMix: Anonymous Messaging via Secure Multiparty Computation. In *Proceedings of the 26th USENIX Conference on Security Symposium*. USENIX Association, 2017.

2. Reyhane Attarian, Esfandiar Mohammadi, Tao Wang, and Emad Heydari Beni. Mixflow: Assessing mixnets anonymity with contrastive architectures and semantic network information. *IACR Cryptol. ePrint Arch.*, page 199, 2023.
3. Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
4. Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
5. Sanjit Bhat, David Lu, Albert Kwon, and Srinivas Devadas. Var-cnn: A data-efficient website fingerprinting attack based on deep learning. *Proceedings on Privacy Enhancing Technologies*, 2019(4):292–310, 2019.
6. Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, 2020.
7. David L. Chaum. Untraceable Electronic Mail, Return Addresses, and Digital Pseudonyms. *Communications of the ACM*, 24, 1981.
8. David L Chaum. Untraceable electronic mail, return addresses, and digital pseudonyms. *Communications of the ACM*, 24(2):84–90, 1981.
9. Chen Chen, Daniele E. Asoni, Barrera, OSC, David, George Danezis, and Adrain Perrig. HORNET: High-Speed Onion Routing at the Network Layer. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. Association for Computing Machinery, 2015.
10. George Danezis. The traffic analysis of continuous-time mixes. In *International Workshop on Privacy Enhancing Technologies*, pages 35–50. Springer, 2004.
11. George Danezis. Traffic Analysis of the HTTP Protocol over TLS, 2009.
12. George Danezis and Ian Goldberg. Sphinx: A compact and provably secure mix format. In *2009 30th IEEE Symposium on Security and Privacy*, pages 269–282. IEEE, 2009.
13. Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
14. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
15. Claudia Diaz, Harry Halpin, Ilya, , and Aggelos Kiayias. The Nym Network: The Next Generation of Privacy Infrastructure. *White Paper, Version 1.0*, 2021.
16. Claudia Diaz and Bart Preneel. Taxonomy of Mixes and Dummy Traffic. In *Information Security Management, Education and Privacy*, 2004.
17. Claudia Díaz, Len Sassaman, and Evelyne Dewitte. Comparison Between Two Practical Mix Designs. In *Computer Security – ESORICS 2004*, 2004.
18. Roger Dingledine and Nick Mathewson. Anonymity loves company: Usability and the network effect. In *WEIS*. Citeseer, 2006.
19. Ran Dubin, Amit Dvir, Ofir Pele, and Ofer Hadar. I know what you saw last minute—encrypted http adaptive video streaming title classification. *IEEE transactions on information forensics and security*, 12(12):3039–3049, 2017.
20. Kristian Gjosteen. The Norwegian Internet Voting Protocol. In *E-Voting and Identity*, 2012.

21. Iness Ben Guirat and Claudia Diaz. Mixnet optimization methods. *Proceedings on Privacy Enhancing Technologies*, 1:22, 2022.
22. Jamie Hayes, George Danezis, et al. k-fingerprinting: A robust scalable website fingerprinting technique. In *USENIX security symposium*, pages 1187–1203, 2016.
23. HOPR. *The Book of HOPR*, 2021. https://hoprnet.org/Book_Of_Hopr_2021.01_v1.pdf.
24. Christiane Kuhn, Martin Beck, Stefan Schiffner, Eduard Jorswieck, and Thorsten Strufe. On Privacy Notions in Anonymous Communication. *Proceedings on Privacy Enhancing Technologies*, 2019.
25. David Lazar, Yossi Gilad, and Nickolai Zeldovich. Karaoke: Distributed private messaging immune to passive traffic analysis. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 711–725, 2018.
26. Xinshu Ma, Florentin Rochet, and Tariq Elahi. Stopping silent sneaks: Defending against malicious mixes with topological engineering. *arXiv preprint arXiv:2206.00592*, 2022.
27. Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An empirical model of large-batch training. *arXiv preprint arXiv:1812.06162*, 2018.
28. Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2013.
29. Brad Miller, Ling Huang, Anthony D Joseph, and J Doug Tygar. I know why you went to the clinic: Risks and realization of https traffic analysis. In *International Symposium on Privacy Enhancing Technologies Symposium*. Springer, 2014.
30. OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
31. Simon Oya, Carmela Troncoso, and Fernando Pérez-González. Do dummies pay off? limits of dummy traffic protection in anonymous communications. In Emiliano De Cristofaro and Steven J. Murdoch, editors, *Privacy Enhancing Technologies*, pages 204–223, Cham, 2014. Springer International Publishing.
32. Simon Oya, Carmela Troncoso, and Fernando Perez-Gonzalez. Understanding the effects of real-world behavior in statistical disclosure attacks. In *2014 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 72–77, 2014.
33. Simon Oya, Carmela Troncoso, and Fernando Pérez-González. Meet the family of statistical disclosure attacks. In *2013 IEEE Global Conference on Signal and Information Processing*, pages 233–236, 2013.
34. Andriy Panchenko, Fabian Lanze, Jan Pennekamp, Thomas Engel, Andreas Zinnen, Martin Henze, and Klaus Wehrle. Website fingerprinting at internet scale. In *NDSS*, 2016.
35. Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014.
36. Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018.
37. Ania M. Piotrowska. Studying the anonymity trilemma with a discrete-event mix network simulator. In *Proceedings of the 20th Workshop on Workshop on Privacy in the Electronic Society, WPES '21*, page 39–44, New York, NY, USA, 2021. Association for Computing Machinery.
38. Ania M. Piotrowska. Studying the Anonymity Trilemma with a Discrete-Event Mix Network Simulator. In *Proceedings of the 20th Workshop on Workshop on Privacy in the Electronic Society*, 2021.

39. Ania M Piotrowska, Jamie Hayes, Tariq Elahi, Sebastian Meiser, and George Danezis. The loopix anonymity system. In *26th USENIX Security Symposium (USENIX Security 17)*, pages 1199–1216, 2017.
40. Fernando Pérez-González, Carmela Troncoso, and Simon Oya. A least squares approach to the static traffic analysis of high-latency anonymous communication systems. *IEEE Transactions on Information Forensics and Security*, 9(9):1341–1355, 2014.
41. Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher. *CoRR*, abs/2112.11446, 2021.
42. Andrei Serjantov and George Danezis. Towards an Information Theoretic Metric for Anonymity. In *Privacy Enhancing Technologies*, 2003.
43. Andrei Serjantov, Roger Dingledine, and Paul Syverson. From a Trickle to a Flood: Active Attacks on Several Mix Types. In *Information Hiding*, 2003.
44. Payap Sirinam, Mohsen Imani, Marc Juarez, and Matthew Wright. Deep fingerprinting: Undermining website fingerprinting defenses with deep learning. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 1928–1943, 2018.
45. Joseph Turian, Lev-Arie Ratinov, and Yoshua Bengio. Word Representations: A Simple and General Method for Semi-Supervised Learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 2010.
46. Jelle Van Den Hooff, David Lazar, Matei Zaharia, and Nickolai Zeldovich. Vuvuzela: Scalable private messaging resistant to traffic analysis. In *Proceedings of the 25th Symposium on Operating Systems Principles*, pages 137–152, 2015.
47. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
48. Tao Wang, Xiang Cai, Rishab Nithyanand, Rob Johnson, and Ian Goldberg. Effective attacks and provable defenses for website fingerprinting. In *Usenix Conference on Security Symposium*, pages 143–157, 2014.
49. Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
50. David Isaac Wolinsky, Henry Corrigan-Gibbs, Bryan Ford, and Aaron Johnson. Scalable anonymous group communication in the anytrust model. In *European Workshop on System Security (EuroSec)*, volume 4, 2012.
51. Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. Harnessing the power of llms in practice: A survey on chatgpt and beyond. 2023.
52. Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A Survey of Large Language Models, 2023.