

Side-channel Attacks and New Principles in the Shuffle Model of Differential Privacy

Shaowei Wang, Jin Li, Changyu Dong, Jin Li, Zhili Zhou, Di Wang, Zikai Wen

Abstract—The shuffle model employs a shuffler to anonymize and permute user messages, thereby enhancing privacy/utility trade-offs compared to the local model. Ideally, it assumes perfect message anonymity protection against adversaries, allowing each user to hide among a large population. However, in contexts like mobile/edge networks or in scenarios where the shuffler is curious, this assumption is frequently unrealistic. In this study, we demonstrate the vulnerability of the shuffle model to communication side-channel attacks, which substantially compromise privacy amplification via shuffling. We categorize side-channel information in the shuffle model into three types: (i) in-out information, revealing the victim user’s participation and timing, (ii) message-cardinality information, indicating the victim’s message count, and (iii) message-length information, disclosing the victim’s message length(s). Numerical results indicate these attacks increase privacy loss by 200% to 4100%, revealing secret value with probability more than 90%. After theoretically analyzing the remaining privacy amplification effects, we suggest several countermeasures and principles to alleviate degradation caused by these attacks: (a) appending padding bits to each message to counter message-length attacks, (b) maximizing query parallelization to elude in-out attacks and increase the population for privacy amplification, and (c) sending dummy messages to exchange communication costs for improved privacy amplification effects. The newly proposed paradigms and principles significantly save privacy budget in comparison to current models under attack.

Index Terms—differential privacy, shuffle model, side channels, privacy attacks

I. INTRODUCTION

THE shuffle model of differential privacy (DP) [1] has emerged as a compelling approach to data privacy protection, combining the benefits of the classical central model [2] (i.e., relatively high data utility) and the local model [3]–[9] (i.e., minimal trust in other parties). In the shuffle model, an intermediary shuffler anonymizes and randomly permutes messages from a user population before forwarding them to the server (e.g., data analysts). As the shuffling operation is data-agnostic and can be executed over ciphertext space, numerous parties can assume the shuffler’s role in practice, including anonymous communication networks [10],

Shaowei Wang, Jin Li, Changyu Dong, Jin Li, and Zhili Zhou are with School of Artificial Intelligence, Guangzhou University. E-mail: Shaowei Wang (wangsw@gzhu.edu.cn).

Di Wang is with King Abdullah University of Science and Technology.

Zikai Wen is with University of Washington (E-mail: zkwen@uw.edu).

This work is supported by National Natural Science Foundation of China (No.62372120, 62472116, 62102108, U23A20307, U21A20463), Guangdong Basic and Applied Basic Research Foundation (No.2023A1515030273, 2022A1515010061), Foundation of Yunnan Key Laboratory of Service Computing (No.YNSC24115), and Science and Technology Projects in Guangzhou (No.2025A03J3182, No.202201010194).

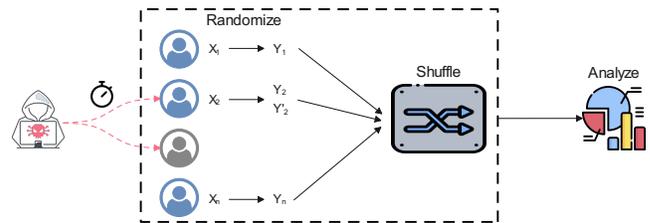


Fig. 1. An illustration of the shuffle model under communication side-channel attacks (the shuffler can be an honest-but-curious adversary).

secure hardware [11], [12], and other cryptographic tools [13]. Moreover, since potential privacy adversaries (e.g., the server) only observe anonymized and shuffled messages, each user can hide within a large population. Consequently, adding a minimal amount of noise to local messages sufficiently protects data privacy in the released view of shuffled messages. This phenomenon is known as *privacy amplification via shuffling* [1]. For example, in the single-message shuffle model where each user sends one message satisfying local ϵ_0 -differential privacy (LDP), the shuffled messages from n users actually preserve $\tilde{O}(\sqrt{e^{\epsilon_0}/n})$ -DP [14]; in the multi-message shuffle model where each user sends multiple messages, as long as these messages are carefully calibrated, the overall shuffled messages can achieve DP. Due to its potential to achieve excellent privacy-utility-efficiency trade-offs in decentralized settings, the shuffle model has been applied across various domains, including count/summation queries [15]–[23] and machine learning [24]–[30], and was deployed in Apple and Google’s Exposure Notification Privacy-preserving Analytics [31].

Despite the success of the shuffle model in decentralized private data analysis, it relies on several unrealistic security assumptions in such environments. Ideally, privacy adversaries in the shuffle model are assumed to lack intermediate information about a victim user during protocol execution, having only access to the crowd’s shuffled messages. We argue that this assumption is easily violated in mobile computing, edge computing, or cable networking environments where the shuffle model is applied (see Figure 1 for an illustration). For example, in wireless/cellular networks, privacy adversaries can precisely infer a victim user’s communication activities at a negligible cost [32]–[36]. Adversaries can effortlessly deduce the victim’s communication timing with the shuffler, potentially ascertain the message payload size, and infer the number of messages contributed by the victim at high precision (e.g., with over 90% accuracy as reported in [33], [34]). *The situation is even worse when the intermediate shuffler is curious: the shuffler observes almost certain side information*

about messages from each user (e.g., message cardinality, approximate message lengths, and the communication timing of a user), even if each message is encrypted by the public key of the server/analyzer.

A. Our Contributions

In this work, we analyze the side-channel vulnerabilities of the shuffle DP model and propose mitigation strategies with formal guarantees. Despite many efforts in the literature on communication side-channel analyses (especially for anonymous channels [37]–[39]), we are the first to provide formal analyses in the differential privacy sense. Additionally, due to differences in privacy goals, we demonstrate that conventional countermeasures (e.g., sending dummy messages [40]) for anonymous networks can induce infinite DP loss or be inefficient in the shuffle model. The contributions of this work are as follows:

Privacy degradation analyses due to side-channel information. We demonstrate that side-channel information, easily inferred by privacy adversaries, significantly impacts privacy amplification effects. Intuitively, a victim user's communication timing reveals the specific task (e.g., the specific gradient descent iteration in federated learning) they participated in, allowing the victim to hide only among the sub-population of that particular task, which comprises only several tens or hundreds of users. If privacy adversaries know the length of the victim's contributed message, they can filter out many other messages with different lengths in the shuffled messages, substantially reducing the number of messages/users the victim can hide among. In certain multi-message protocols, the number of messages sent by a user can be directly related to the user's true value. Consequently, the knowledge of message cardinality by privacy adversaries completely undermines privacy amplification through shuffling. Numerical results on typical settings show that communication timing (termed as in-out attacks) and message-length attacks increase the victim user's differential privacy loss by 200% to 4100%, revealing the victim's secret value with probability more than 90%; for several state-of-the-art multi-message protocols, the message-cardinality attack increases privacy loss to $+\infty$, revealing the victim's secret value with probability more than 95%.

Defending against side-channel attacks. We conducted formal analyses of privacy degradation in widely-used shuffle models (e.g., shuffle-then-randomize model [1], [14], [41] and randomize-then-shuffle model [15], [24], [43]) under side-channel attacks and developed new shuffle models resistant to such attacks. In the proposed model, users make their own participation choices, reducing trust assumptions regarding the shuffler/server. To minimize privacy loss due to participation timing, the original data analysis task is first represented as a directed acyclic graph (DAG). Based on this, multiple independent queries are issued simultaneously, while sequential dependent queries are addressed using dummy messages from users, ensuring indistinguishable participation timings. To counteract privacy loss resulting from message payload size/length, each user's message is padded to a fixed length (as the traditional defensive methods in the literature [44] for

anonymous channels). To address severe privacy leaks from message cardinality where traditional defenses fail, domain transformation can be employed specific to the multi-message protocol. These techniques form new principles for the shuffle models under side-channel attacks. Furthermore, even in settings without attacks, our proposed multinomial-randomize-shuffle model provides much stronger privacy amplification effects than existing models.

We summarize the results of in-out/message-length attacks and defenses in Table I. The amplification population size denotes the number of users a victim can hide among - a critical parameter in privacy amplification via shuffling, where higher values indicate better amplification effects. The variables include: n , the full population size; n_k , the size of the population in the same k -th batch as the victim user; $n_{*,l}$, the size of the population with the same message length as the victim user; $n_{k,l}$, the size of the population with the same message length and in the same k -th batch as the victim user; and $n_{(m)}$, the size of the population that selected the same m -th bin as the victim user. The *normal* column denotes the amplification population size when there are no side-channel attacks; the *in-out attacks* column represents the amplification size under in-out attacks; the *in-out & length atk.* column signifies the amplification size under joint in-out and message-length attacks. The *adaptive query* column indicates whether the model supports sequentially adaptive queries. The *extra costs* column highlights the additional communication costs compared to the normal model.

B. Organization

The remainder of this paper is organized as follows. Section II reviews related work on the shuffle model and side channels. Section III provides background knowledge. Section IV identifies several side-channel attacks and numerically analyzes their impact on privacy. Section V presents theoretical analyses of the privacy degradation impact of side-channel information. Section VI proposes defense methods. Section VII recaps new requirements for the shuffle model, and discusses the implication for the broader DP-Crypto systems. Finally, Section VIII concludes the paper.

II. RELATED WORK

A. Side-channel Attacks and Defenses

Side-channel attacks exploit unintended information leakage through a system's physical properties or observable characteristics, such as electromagnetic radiation, power consumption, and communication patterns.

Side channels of anonymous networks. Anonymous networks [10] aim to protect users' identities and relationships between communicating parties. Researchers have investigated various side-channel attacks that target anonymous channels, including traffic analysis attacks [37], timing attacks [38], and intersection attacks [39]. These attacks exploit information leaks from communication patterns, message timings, or user behavior to undermine the anonymity guarantees provided by the channels. As countermeasures, researchers propose strategies like padding messages, traffic camouflage [44], sending

TABLE I
COMPARISON OF THE AMPLIFICATION POPULATION SIZE AND OTHER PROPERTIES OF VARIOUS SHUFFLE MODELS UNDER IN-OUT AND MESSAGE-LENGTH ATTACKS. CONVENTIONALLY, $n \geq n_{(m)} \gg n_k \geq n_{k,l} \gg 1$, INDICATING DISPARATE AMPLIFICATION EFFECTS OF VARIOUS MODELS UNDER ATTACKS.

shuffle models	normal setting	in-out attacks	in-out & length atk.	adaptive queries	extra costs
shuffle-then-randomize [1], [14], [41][Definition 5]	n	1	1	✓	N.A.
divide-randomize-shuffle [25]–[30] [Definition 6]	n_k	n_k	$n_{k,l}$	✓	N.A.
subsample-randomize-shuffle [24], [42] [Definition 7]	n	n_k	$n_{k,l}$	✓	N.A.
parallel-randomize-shuffle [1], [43][Algorithm 3]	n	n	$n_{*,l}$	✗	N.A.
multinomial-randomize-shuffle this work [Definition 11]	n	n_k	$n_{k,l}$	✓	N.A.
rectified parallel-randomize-shuffle this work [Algorithm 4]	n	n	n	✗	padding bits per message
rectified multinomial-randomize-shuffle this work [Algorithm 5]	n	n	n	✓	padding bits per message K -time dummy messages
rectified bin-randomize-shuffle this work [Algorithm 6]	$n_{(m)}$	$n_{(m)}$	$n_{(m)}$	✓	padding bits per message few dummy messages

dummy messages [40] and disturbing sending times [45] in anonymous networks.

This work presents the first side-channel attack/defense study to joint systems of cryptographic tools (i.e., the shuffling/anonymous network) and information-theoretic privacy tools. We provide formal privacy analyses of side-channel attacks in the shuffle model within the DP context. We also propose defense methods incorporating both classical countermeasures (e.g., padding messages and sending dummy messages [40]) and tailored approaches. Specifically, due to the differences in privacy goals, we demonstrate that certain classical countermeasures for anonymous networks can induce infinity DP loss (see Section VI-D) or being inefficient (see Section VI-C) in the shuffle model.

B. Shuffle Model of Differential Privacy

The shuffle model of DP anonymizes user messages through an intermediate shuffler before sending them to a server for analytics. The model’s foundation lies in privacy amplification analysis [15], [41], [43], ensuring the global privacy level of shuffled messages. Depending on the number of messages a user can send, the shuffle model can be categorized as either multi-message [21], [46] or single-message [1], [15], [41].

Various shuffle model variants exist in the literature. For convenient theoretical analysis of privacy amplification, studies such as [1], [14], [41] utilize the shuffle-then-randomize model, in which user data are first shuffled and then fed into adaptive local randomizers sequentially. Practical approaches to support multiple adaptive queries in decentralized setting is to divide the user population into non-overlapping parts (e.g., in [25], [27], referred to as the divide-randomize-shuffle model), or letting the shuffler to perform user subsampling (e.g., in [24], [42], referred to as the subsample-randomize-shuffle model). Another approach is to let each user randomly select (at most) one query to check in [16], and then choose one user (from checked-in users) for each query slot.

We demonstrate that aforementioned shuffle models are all susceptible to side-channel attacks. Their privacy amplification

degrades to the intra-batch level or even to the local level (i.e., no privacy amplification), refer to Table I for summarized results. Worse still, for several protocols in these shuffle models, such as SOTA multi-message summation protocols [17], [19], [46], [47], the differential privacy loss under side-channel attack increases to $+\infty$.

C. Security Attacks of Differential Privacy

Differential privacy, as the de facto standard for data privacy, is widely adopted in industry for sensitive databases [48] and decentralized data collection/analyses [49]. Although DP provides rigorous data privacy in theory, its practical implementation may encounter unexpected security issues. For example, improper implementations of floating-point numbers for the prevalent Laplace/Gaussian mechanism compromise the intended DP [50], [51]. There can also be side-channel timing/global state attacks in central DP database engines [52]. Besides, a small portion of adversarial users may completely undermine the aggregation results of locally private protocols [53] and shuffle private protocols [18]. There also might be collusion between the shuffler and other parties [54]. To the best of our knowledge, this work is the first to study side-channel security issues of decentralized DP, and challenges the honest-and-not-curious assumption on the shuffler in the shuffle model. As comparison, in the local model of DP [3], where each user trusts no other party, there is limited vulnerability to side-channel attacks, as the plaintext data leaving the user/device is publishable, and thus side-channel information related to it is also publishable.

III. PRELIMINARIES

This section presents definitions of differential privacy and shuffle models. Notations are listed in Table II.

A. Divergences and Differential Privacy

Definition 1 (Hockey-stick divergence): The Hockey-stick divergence between two random variables P and Q is:

$$D_{e^\epsilon}(P||Q) = \int \max\{0, P(x) - e^\epsilon Q(x)\} dx,$$

TABLE II
LIST OF NOTATIONS.

Notation	Description
$[i]$	$\{1, 2, \dots, i\}$
$[i : j]$	$\{i, i + 1, \dots, j\}$
n	the number of users (data owners)
K	the number of sequential queries
U_k	the users participated in the k -th round query
\mathbb{X}	the domain of input data
\mathbb{Y}	the domain of a single message
\mathbb{Z}	the domain of an analyzing algorithm
ϵ_0	the local privacy budget
δ	the failure/violating probability
\mathcal{S}	the shuffling procedure
\mathcal{R}	the randomization algorithm
\mathcal{A}	the data analyzing algorithm
io_v	the participation information of user v
num_v	the number of messages from user v
$len(y)$	the message length (bits) of message y

where P and Q denote both the random variables and their probability density functions.

Two variables P and Q are (ϵ, δ) -indistinguishable if $\max\{D_{e^\epsilon}(P||Q), D_{e^\epsilon}(Q||P)\} \leq \delta$. For two datasets of equal size that differ only by a single individual's data, they are referred to as *neighboring datasets*. Differential privacy limits the divergence of query results on neighboring datasets (see Definition 2). Similarly, in the local setting that accepts a single individual's data as input, we introduce the local (ϵ, δ) -differential privacy in Definition 3. When $\delta = 0$, the concept is abbreviated as ϵ -LDP.

Definition 2 (Differential privacy [2]): A protocol $\mathcal{R} : \mathbb{X}^n \mapsto \mathbb{Z}$ satisfies (ϵ, δ) -differential privacy if, for all neighboring datasets $X, X' \in \mathbb{X}^n$, $\mathcal{R}(X)$ and $\mathcal{R}(X')$ are (ϵ, δ) -indistinguishable.

Definition 3 (Local differential privacy [3]): A protocol $\mathcal{R} : \mathbb{X} \mapsto \mathbb{Y}$ satisfies local (ϵ, δ) -differential privacy if, for all $x, x' \in \mathbb{X}$, $\mathcal{R}(x)$ and $\mathcal{R}(x')$ are (ϵ, δ) -indistinguishable.

B. The Classic Shuffle Model

This part reviews the classic one-round shuffle model.

Single-message shuffle model. Following the conventions of the randomize-then-shuffle model [15], [55], we define a single-message protocol \mathcal{P} as a list of algorithms $\mathcal{P} = (\{\mathcal{R}_i\}_{i \in [n]}, \mathcal{A})$, where $\mathcal{R}_i : \mathbb{X} \rightarrow \mathbb{Y}$ is user i 's local randomizer, and $\mathcal{A} : \mathbb{Y}^n \rightarrow \mathbb{Z}$ is the analyzer on the data collector's side. The overall protocol implements a mechanism $\mathcal{P} : \mathbb{X}^n \rightarrow \mathbb{Z}$ as follows. User i holds a data record x_i and a local randomizer \mathcal{R}_i , then computes a message $y_i = \mathcal{R}_i(x_i)$. The messages y_1, \dots, y_n are shuffled and submitted to the analyzer. We denote the shuffling step as $\mathcal{S}(y_1, \dots, y_n)$, where $\mathcal{S} : \mathbb{Y}^n \rightarrow \mathbb{Y}^n$ is a *shuffler* that applies a uniform-random permutation to its inputs. In summary, the output of $\mathcal{P}(x_1, \dots, x_n)$ is represented by $\mathcal{A} \circ \mathcal{S} \circ \mathcal{R}_{[n]}(X) = \mathcal{A}(\mathcal{S}(\mathcal{R}_1(x_1), \dots, \mathcal{R}_n(x_n)))$.

Multi-message shuffle model. In contrast to sending a single message, the multi-message shuffle model allows each user to release multiple messages to the shuffler. The $\mathcal{R}_i : \mathbb{X} \rightarrow \mathbb{Y}^*$ is user i 's local randomizer. The output $\mathcal{P}(x_1, \dots, x_n)$ of the

overall protocol is $\mathcal{A} \circ \mathcal{S} \circ \mathcal{R}_{[n]}(X) = \mathcal{A}(\mathcal{S}(\mathcal{R}_1(x_1) \cup \mathcal{R}_2(x_2) \cup \dots \cup \mathcal{R}_n(x_n)))$.

The shuffle model strives to ensure the privacy of $\mathcal{P}(x_1, \dots, x_n)$ for any analyzer \mathcal{A} . Owing to the post-processing property of Hockey-stick divergence, guaranteeing that the shuffled messages $\mathcal{S} \circ \mathcal{R}_{[n]}(X)$ exhibit differential privacy suffices. We formally delineate differential privacy in the shuffle model in Definition 4.

Definition 4 (DP in the shuffle model): A protocol $\mathcal{P} = (\{\mathcal{R}_i\}_{i \in [n]}, \mathcal{A})$ satisfies (ϵ, δ) -DP in the shuffle model iff for all neighboring datasets X and $X' \in \mathbb{X}^n$, the $\mathcal{S} \circ \mathcal{R}_{[n]}(X)$ and $\mathcal{S} \circ \mathcal{R}_{[n]}(X')$ are (ϵ, δ) -indistinguishable.

C. Prevalent Variants of Shuffle Model

Several variants of shuffle model exist in the literature, depending on the organization of users to respond to multiple adaptive queries.

1) Shuffle-then-randomize Model: We revisit the ideal (single-message) shuffle model based on shuffle-then-randomize [1], [14], [41]. Given an input dataset $X = \{x_1, \dots, x_n\}$, a uniform-random permutation $\pi : [n] \mapsto [n]$ is first applied to obtain $\mathcal{S}(X) = \{x_{\pi^{-1}(1)}, \dots, x_{\pi^{-1}(n)}\}$, followed by a series of adaptive randomizers $\{\mathcal{R}_i\}_{i \in [n]}$. The i -th randomizer \mathcal{R}_i takes the i -th datum $x_{\pi^{-1}(i)}$ in $\mathcal{S}(X)$ and the previous $i - 1$ randomization results as input (see Definition 5). Since π is not revealed to the server, each user's message can hide among all n messages. Specifically, when every \mathcal{R}_i satisfies ϵ_0 -LDP (not necessarily identical), the messages $\mathcal{R}_{[n]} \circ \mathcal{S}(X)$ satisfy $(\tilde{O}(e^{\epsilon_0/2}/\sqrt{n}), \delta)$ -DP [14], [41]. This model serves as an ideal one that supports fully adaptive queries while providing the strongest privacy amplification effects (amplified by n users).

Definition 5 (Shuffle-then-randomize model [1], [14], [41]): Let $\mathcal{R}_i : \mathbb{Z}_0 \times \mathbb{Y}_1 \times \dots \times \mathbb{Y}_{i-1} \times \mathbb{X} \rightarrow \mathbb{Y}_i$ for $i \in [n]$ be a sequence of algorithms, where \mathbb{Z}_0 denotes the range space of global information. A protocol $\mathcal{P}_{s-r} : \mathbb{Z}_0 \times \mathbb{X}^n \rightarrow \mathbb{Y}_0 \times \dots \times \mathbb{Y}_n$ with global information z_0 in the shuffle-then-randomize model proceeds as follows: given a dataset $x_{[n]} \in \mathbb{X}^n$, it samples a uniform-random permutation π and then sequentially computes $z_i = \mathcal{R}_i(z_{[0:i-1]}, x_{\pi^{-1}(i)})$ for $i \in [n]$ before outputting $z_{[0:n]}$.

2) Divide-randomize-shuffle Model: In decentralized settings, a more realistic approach is letting the analyzer divide users into multiple groups and employ each group for one query sequentially, see Definition 6. This model is adopted by many works for federated learning [27], and for bandit/reinforcement learning [25], [28]–[30]. Since the division is known to the analyzer (i.e., a potential adversary), a user $i \in U_k$ can only hide among users in the same group, thus privacy can only be amplified by $|U_k|$. This is considerably weaker than the ideal shuffle model.

Definition 6 (Divide-randomize-shuffle model [25], [26]): Let $U_{[K]}$ denote a division of set $[n]$ such that $U_k \cap U_{k'} = \emptyset$ for all $k, k' \in [K]$ that $k \neq k'$, and $U_1 \cup \dots \cup U_K = [n]$. Let $\mathcal{R}_{(k)} : \mathbb{Z}_0 \times \mathbb{Z}_1 \times \dots \times \mathbb{Z}_{k-1} \times \mathbb{X} \rightarrow \mathbb{Y}_{(k)}$ denote the randomizer in the k -th round, where $\mathbb{Z}_k = \mathbb{Y}_{(k)}^{|U_k|}$ and $\mathbb{Y}_{(k)}$ is the range space of $\mathcal{R}_{(k)}$. The $z_{(0)} \in \mathbb{Z}_0$ is global information. A protocol

$\mathcal{P}_{d-r-s} : \mathbb{Z}_0 \times \mathbb{X}^n \rightarrow \mathbb{Z}_0 \times \cdots \times \mathbb{Z}_K$ in the divide-randomize-shuffle model operates as follows: given a dataset $x_{[n]} \in \mathbb{X}^n$, it sequentially computes $z_{(i)} = \mathcal{S}(\mathcal{R}_{(k)}(z_{(0:k-1)}, x_i)_{i \in U_k})$ and outputs $z_{(0)}, z_{(1)}, \dots, z_{(K)}$ along with $U_{[K]}$.

3) *Subsample-randomize-shuffle Model*: Rather than letting the analyzer (i.e., a potential adversary) to perform user division, a more practical method involves having the shuffler randomly sample users in each round independently, as illustrated in Definition 7. This approach is commonly employed for federated learning, as seen in [24], [42].

Definition 7 (*Subsample-randomize-shuffle model* [24], [42]): Let $\{U_k\}_{k \in [K]}$ denote a list of subsets that each has size s and is uniform-randomly sampled from $[n]$ without replacement. Let $\mathcal{R}_{(k)} : \mathbb{Z}_0 \times \mathbb{Z}_1 \times \cdots \times \mathbb{Z}_{k-1} \times \mathbb{X} \rightarrow \mathbb{Y}_{(k)}$ denote the randomizer in the k -th, where $\mathbb{Z}_k = \mathbb{Y}_{(k)}^{U_k}$ and $\mathbb{Y}_{(k)}$ is the range space of $\mathcal{R}_{(k)}$. The $z_{(0)} \in \mathbb{Z}_0$ is global information. A protocol $\mathcal{P}_{d-r-s} : \mathbb{Z}_0 \times \mathbb{X}^n \rightarrow \mathbb{Z}_0 \times \cdots \times \mathbb{Z}_K$ in the subsample-randomize-shuffle model operates as follows: given a dataset $x_{[n]} \in \mathbb{X}^n$, it independently samples $\{U_k\}_{k \in [K]}$ as described previously, then sequentially computes $z_{(i)} = \mathcal{S}(\mathcal{R}_{(k)}(z_{(0:k-1)}, x_i)_{i \in U_k})$ and outputs $z_{(0)}, \dots, z_{(K)}$.

Given that only a relatively small batch of users, with a size of s , is randomly selected for each round, privacy is also amplified by subsampling [56]. Assuming that local randomizers satisfy ϵ_0 -LDP, and when combined with privacy amplification via shuffling and the advanced composition theorem of differential privacy, the overall privacy loss is $(\tilde{O}(\frac{s}{n} \sqrt{K e^{\epsilon_0}/s}), \delta)$ [24], [42].

IV. SIDE-CHANNEL ATTACKS IN THE SHUFFLE MODEL

In this part, we show privacy amplification is severely damaged by side-channel knowledge, which can be readily divulged to privacy adversaries in decentralized settings.

A. Threat model of side-channel attacks

We assume that the internal shuffling process \mathcal{S} is perfectly secure, with adversaries only able to observe the victim users' communication activities (with the shuffler) and the output $\mathcal{P}(X)$ of the shuffle model. Typically, every user in the shuffle model encrypts their messages $\mathcal{R}_i(x_i)$ using the analyzer's public key before sending them to the shuffler. Privacy adversaries (including a potential adversary, the curious shuffler) do not have access to the values of these encrypted messages, but only to their side information, such as communication timing, the number of messages, and message length(s).

The aim of privacy adversaries is to deduce the secret x_i of victim users, utilizing either the success probability of inferring x_i or the differential privacy loss of x_i as a metric.

B. Categories of Side-channel Information

In the shuffle models described above, we define the following three types of side-channel information about the victim user v ($v \in [n]$):

I. (*In-out information*) In-out information indicates the victim's participation in multiple rounds of a shuffle protocol. We denote it as $io_v \in \{0, 1\}^K$, where the k -th value

$io_v(k)$ is 1 if the victim participated in the k -th round, and 0 otherwise.

II. (*Message-cardinality information*) In a multi-message protocol, message-cardinality information indicates the number of messages from the victim. We denote this as $num(\mathcal{R}_v(x_v))$, representing the number of messages output by $\mathcal{R}_v(x_v)$.

III. (*Message-length information*) Message-length information indicates the length of the victim's message. We denote this as $len(\mathcal{R}_v(x_v))$, representing the number of bits in $\mathcal{R}_v(x_v)$. In a multi-message protocol, this represents the length of each message from $\mathcal{R}_v(x_v)$.

We argue that in decentralized settings, such as mobile computing and wireless networks, privacy adversaries can easily infer in-out, message-cardinality, and message-length information. The in-out information is strongly correlated with the communication patterns of the victim user, which are almost public information in wireless networks or cellular networks [57], [58]. Adversaries can sniff the number of packets sent from the victim to the shuffler through network activities and by examining packet headers, even when packets are delivered over the prevalent HTTPS [59]. Furthermore, the length of TCP/IP packets from the victim can be monitored by sniffing network traffic [60]–[62]. Relatively, the in-out information is easier to be inferred. The side-channel information pertaining to message number or length might often be imprecise. This lack of accuracy is predominantly due to the implementation of techniques such as packet padding and packet slicing by modern wireless or cellular networking protocols. Additionally, this imprecision is posed by modern encryption schemes, which encrypt payloads on a block-by-block basis. Notably, when the shuffler is a curious party, since the shuffler observes encrypted message(s) from each individual, the leaked side-channel information can be quite precise, regardless of the communication media between each user and the shuffler.

Note that in-out information is especially vulnerable in interactive queries with substantial download overheads, such as federated machine learning. Users download the latest model parameters only when participating in a round. The correlated heavy-loaded downloading traffic, along with uploading traffic to the shuffler, significantly increases the risk of in-out information exposure.

C. Privacy under Message-cardinality Attacks

Message-cardinality attacks occur in multi-message protocols, where the number of messages a user sends to the shuffler might depend on the actual value the user holds (e.g., in [17], [19], [27], [46], [47], [63]). If message cardinality is observed by adversaries, privacy can be compromised.

As an example, we show a message-cardinality attack to the state-of-the-art Δ -summation protocol [19]. In this protocol, each user holds an integer $x_i \in [0, \Delta]$, and the analyzer obtains the sum satisfying (ϵ, δ) -DP. The algorithms of the local randomizer and the analyzer are shown in Algorithm 1 and 2. The local randomizer randomizes the user's input in three steps: (1) each user i sends x_i if it is non-zero (line 1-2); (2)

it sends unary noise messages (i.e. z_i^{+1} copies of +1 and z_i^{-1} copies of -1 following negative binomial distributions) whose sum is equal to the Discrete Laplace noise commonly used algorithms in the central DP model (line 3-4); (3) it defines a sub-collection S of all multisets of $\{-\Delta, \dots, +\Delta\} \setminus \{0\}$ whose sum of elements is equal to zero (e.g. $\{-1, -1, 2\}$), and for each multiset $s \in S$, z_i^s is sampled and z_i^s copies of every element in s are sent (line 5-8). The analyzer simply sums up all messages.

Algorithm 1: Δ -Summation Randomizer [19]

```

1 if  $x_i \neq 0$  then
2   Send  $x_i$ 
3 Sample  $z_i^{+1}, z_i^{-1} \sim \text{NB}(1, e^{-(1-\gamma)\epsilon})/n$ 
4 Send  $z_i^{+1}$  copies of +1, and  $z_i^{-1}$  copies of -1
5 for  $s \in S$  do
6   Sample
7   for  $m \in s$  do
8     Send  $z_i^s$  copies of  $m$ 

```

Algorithm 2: Δ -Summation Analyzer [19]

```

1  $T \leftarrow$  multiset of messages received
2 return  $\sum_{y \in T} y$ 

```

For simplicity, let us consider the case $\Delta = 1$. Now the zero-sum messages collection S simply contains one message set, i.e. $S = \{\{-1, 1\}\}$. Let $Z^1 \sim \text{NB}(1, e^{-(1-\gamma)\epsilon})/n$, $Z^2 \sim \text{NB}(1, e^{-(1-\gamma)\epsilon})/n$, and $Z^3 \sim \text{NB}(3(1 + \log(2/\delta)), e^{-0.1 \min(1, \gamma\epsilon)/4})/n$, then when $x_i = 0$, the total number of messages is $\text{num}_0 = Z^1 + Z^2 + 2Z^3$ where $Z^1 + Z^3$ is the number of +1, and $Z^2 + Z^3$ is the number of -1. On the other hand, when $x_i = 1$, the total number of messages is $\text{num}_1 = 1 + Z^1 + Z^2 + 2Z^3$ because the local randomize also sends x_i in addition to the noise messages. Our observation is that num_0 can be 0 and num_1 cannot, and in typical settings, $\text{num}_b = b$ with a large probability. In short, the number of messages exposes the user's private data value. To demonstrate this, in Figure 2(a), we plot the probability distributions of $\text{num}_0, \text{num}_1$ under the following parameters: $n = 10^4$, $\epsilon = 1.0$, $\delta = 0.01/n$, and $\gamma = 0.1$. As we can see in the figure, $\mathbb{P}[\text{num}_b = b] \geq 95.1\%$ for $b \in \{0, 1\}$. So if an adversary can observe the message cardinality, it can deduce almost certainly the user's private data value. A similar conclusion holds for multi-message protocols in [17], [27], [46], [47], [63].

D. Privacy under In-out Attacks

This section focuses on the privacy degradation of shuffle models under in-out attacks. We start with the ideal shuffle model: shuffle-then-randomize, and then give results about the prevalent subsample-randomize-shuffle model.

Shuffle-then-randomize model under in-out attacks [no privacy amplification]. The shuffle-then-randomize model,

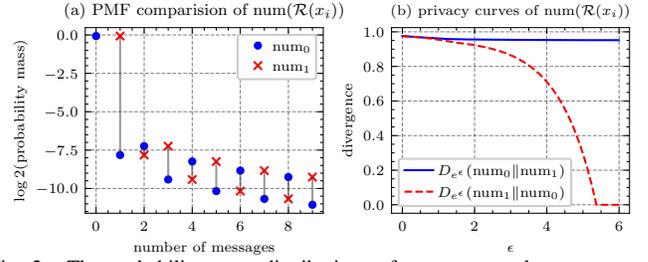


Fig. 2. The probability mass distributions of message numbers $\text{num}_0 = \text{num}(\mathcal{R}(0))$, $\text{num}_1 = \text{num}(\mathcal{R}(1))$ and privacy curves due to message-cardinality attack on [19].

introduced by the seminal work [1] and refined by recent studies [14], [41], provides the best results to date in terms of privacy amplification. In this model, there are multiple rounds. In each round, the analyzer sends a query to the shuffler, who then chooses a user to answer the query according to a random permutation π . The i -th user answers the query in the $\pi(i)$ -th round by sending the locally randomized data to the shuffler. The query and the local randomizer may vary in each round.

In this model, privacy amplification can be achieved because the users can hide themselves in a crowd while answering queries. However, if the adversary can observe a user participating in a particular round (i.e. an in-out attack), privacy is then reduced to whatever the local randomizer can provide and amplification is no longer achievable. We can demonstrate the difference through an example. Consider a scenario in which an analyzer queries a crowd of $n = 10,000$ users, and each of them holds a bit. In each round, the analyzer asks a user to submit her bit through the shuffler. Each user can use a local randomized response mechanism to perturb her bit before submitting it to the shuffler:

$$\mathcal{R}_i(x_i) = \begin{cases} x_i, & \text{with probability } 1 - p; \\ 1 - x_i, & \text{with probability } p. \end{cases}$$

where $p = \frac{1}{e^{\epsilon_0} + 1}$ and ϵ_0 is the local privacy budget. Suppose the global privacy goal is $(\epsilon = 0.2, \delta = 10^{-6})$ -DP, then the local level only needs to be $\epsilon_0 = 2.81$ based on the analysis in [14]. An in-out attack means that the adversary can recover the user's data much easier: the i -th user submitted her data in the $\pi(i)$ -th round, which is observed by the adversary. Since the local randomizer's budget is $\epsilon_0 = 2.81$, the probability $p \approx 0.057$. Hence, the adversary knows that with a high probability $1 - p = 0.943$, the i -th user's bit is the one submitted in the $\pi(i)$ -th round.

Similar attacks apply to binary randomized response used in [1] for streaming data aggregation, generalized randomized response used by [14], [41] for histogram estimation, or any other local randomizers in the shuffle-then-randomize model.

Subsample-randomize-shuffle model under in-out attacks [intra-batch amplification]. In the subsample-randomize-shuffle model, a randomly sampled sub-population report their data in each round. Privacy amplification in this model relies on the fact that the analyzer does not know who participated in which round, hence can amplify privacy also by subsampling [56]. However, if an adversary possesses the victim user's in-out information, the victim's privacy will be weakened significantly. Assuming the victim participated in

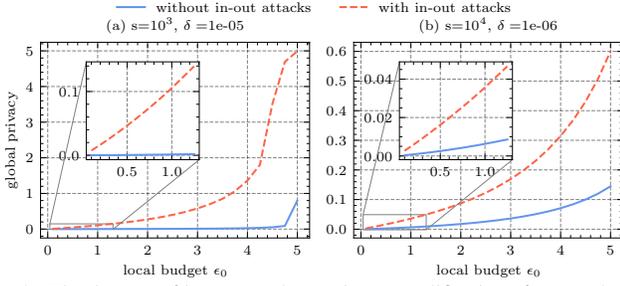


Fig. 3. The damage of in-out attack on privacy amplification of protocol [24] in subsample-randomize-shuffle model.

the k_v -th round, although the adversary cannot identify the exact message sent by the victim, due to the fact that U_{k_v} is a smaller sub-population and is known by the adversary, privacy amplification via subsampling is no longer available.

To demonstrate, let us consider an example from [24]. In this example, subsample-randomize-shuffle is used for training neural networks through federated learning. To jointly train a classifier using MNIST dataset, each user holds a gradient vector $x_i \in \{-0.01, 0.01\}^d$, where $d = 13170$. In each round, $s = 1000$ users are selected from a total $n = 60000$ population. They sanitize their gradient vectors using randomized response on a randomly selected dimension with budget $\epsilon_0 = 2$. Without the in-out attack, a single round gradient aggregation consumes a DP budget (i.e. incurs a privacy loss) of $(0.0064, 10^{-5})$ (using the near-optimal shuffle amplification upper bound in [14, Theorem 3.1] and subsampling amplification bound in [24, Lemma 3]). However, when the adversary possesses the in-out information of a victim participated in a round, the actual privacy budget increases by over 41 times to $(0.27, 10^{-5})$ (using the shuffle amplification lower bound for $s = 1000$ users with randomized response in [43, Theorem 5.1]). Figure 3 shows more privacy loss comparisons with and without in-out attacks. In the worst-case, if a user v is selected for $\Omega(K)$ rounds, the accumulated privacy loss is $\Omega(\sqrt{K}e^{\epsilon_0}/s)$, creating a gap of s/n from the expected privacy loss without in-out attack.

E. Privacy under Message-length Attacks

This section studies the impact of message-length leakage on privacy in the shuffle model, analyzing both local privacy and shuffle privacy amplification impacts.

Privacy amplification degradation. Intuitively, if privacy adversaries know the victim's message length $\mathcal{R}_i(x_i)$, they can filter out messages in $\mathcal{S} \circ \mathcal{R}_{[n]}(X)$ with unmatched lengths. Thus, the victim's message can only hide among a subset users with the same length. Specifically, message length information reduces the population size for privacy amplification to:

$$\#\{\mathcal{R}_i(x_i) \mid i \in [n_k] \setminus \{v\} \text{ and } \text{len}(\mathcal{R}_i(x_i)) = \text{len}(\mathcal{R}_v(x_v))\}. \quad (1)$$

For instance, the seminal work [1] on shuffle privacy amplification considers longitude data $x_i \in \{0, 1\}^d$ aggregation over a period of time $[1 : d]$. To avoid $\Theta(d)$ errors in estimators, a common practice is representing x_i by its hierarchical residue and having users report one hierarchy level. Assuming periods time as $d = 2^H$, the k -th value in the h -th hierarchy is

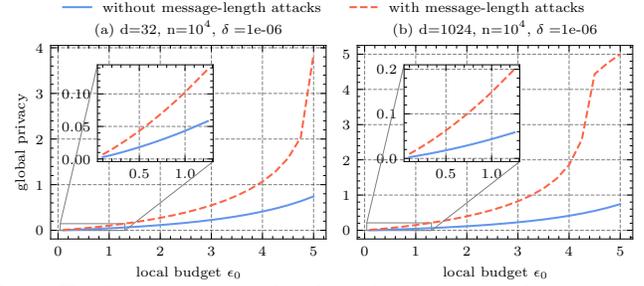


Fig. 4. The damage of message length attack on privacy amplification of the single-message protocol [1].

given by $V_{h,k} = x_i(k \cdot 2^h) - x_i((k-1) \cdot 2^h)$ (assumed $x_i(0) = 0$), where $h \in [0 : H]$, $k \in [1 : d/2^h]$. The study lets users uniformly select one hierarchy level $h \in [0 : H]$ and employ binary randomized response with full ϵ_0 to sanitize ternary vector $V_{h,*} \in \{-1, 0, 1\}^{d/2^h}$, obtaining a message $Y_{h,*} \in \{-1, 1\}^{d/2^h}$. Depending on the hierarchy level h selected, the message length of $Y_{h,*}$ varies significantly, ranging from 1 to d . Consequently, the victim can only hide among approximately $n/(\log_2 d + 1)$ users. In Figure 4, we compare remaining privacy protection under message-length attacks with claimed privacy levels in [1]. The actual privacy loss increases by over 200%, and the gap grows with d . Similar issues exist in other prevalent aggregation tasks (e.g., for range queries [64] and marginal queries [65]) employing sampling-based query selection or compressed binary randomized response [3]. In these protocols, while $\text{len}(\mathcal{R}_i(x_i))$ follows the same distribution P_{len} for any input data x_i , message length may vary widely with high entropy. When adversaries possess the message length $\text{len}(\mathcal{R}_v(x_v)) = l$ of the victim user v , the user can only hide among a much smaller randomized sub-population (of size $|U_{*,l}| \approx 1 + (n-1) \cdot P_{\text{len}}[l]$).

Local privacy loss. When local randomizers' outputs are compressed (e.g., via list representation of sparse ϵ_0 vector in RAPPOR [49], key-value data randomizer [66]) for transmission efficiency as in [27], the compressed message length might probabilistically reveal information about the message value, thus incur severe local privacy loss (see Appendix A).

Joint in-out and message-length Attack. The timing of message transmission (in-out information) and the size of transmitted messages (message-length information) are often simultaneously leaked to adversaries. Consequently, the population a victim can hide among is further restricted to those sharing the same message length in the same round, and the privacy amplification effect under joint in-out and message-length attacks can be further diminished.

V. THEORETICAL ANALYSES OF SIDE-CHANNEL ATTACKS

This section presents formal analyses of the detrimental impact of side-channel information on privacy guarantees in shuffle models. To facilitate theoretical characterization, we introduce some properties of distance measures employed in DP. The data processing inequality asserts that the privacy guarantee cannot be weakened by further analysis of a privatization mechanism's output.

Definition 8 (Data processing inequality [67]): A distance measure $D : \Delta(\mathbb{T}) \times \Delta(\mathbb{T}) \rightarrow [0, \infty]$ on the probability

distribution space satisfies data processing inequality if, for all distributions P and Q in $\Delta(\mathbb{T})$ and for all (possibly randomized) functions $g : \mathbb{T} \rightarrow \mathbb{T}'$,

$$D(g(P)\|g(Q)) \leq D(P\|Q).$$

In the shuffle model, a plethora of sources of randomness exists (e.g., randomness of users' participation choices and randomness of query results from previous rounds). We put forth two instrumental tools for scrutinizing the distance measures under intricate sources of randomness: separability property (refer to Definition 9), and conditioning increasing property (refer to Definition 10).

Definition 9 (Separability property [67]): A distance measure $D : \Delta(\mathbb{T}) \times \Delta(\mathbb{T}) \rightarrow [0, \infty]$ on the space of probability distributions satisfies separability property if, for all distributions P and Q that are joint densities over $\mathbb{T} = \mathbb{T}_1 \times \mathbb{T}_2$ with the same marginal density with respect to \mathbb{T}_1 , i.e. $P = P_{T_1} \cdot P_{T_2|T_1}$ and $Q = P_{T_1} \cdot Q_{T_2|T_1}$,

$$D(P\|Q) = \mathbb{E}_{t_1 \sim P_{T_1}} [D(P|t_1\|Q|t_1)],$$

where $P|t_1$ and $Q|t_1$ denote the conditional variables $P_{T_2|t_1}$ and $Q_{T_2|t_1}$, respectively.

Definition 10 (Conditioning increasing property [67]): A distance measure $D : \Delta(\mathbb{T}) \times \Delta(\mathbb{T}) \rightarrow [0, \infty]$ satisfies condition increase property if, for all distributions P and Q over \mathbb{T} that are generated by $P = \int P_{T|T_1} dT_1$ and $Q = \int Q_{T|T_1} dT_1$ with the same variable T_1 ,

$$D(P\|Q) \leq \mathbb{E}_{t_1 \sim P_{T_1}} [D(P|t_1\|Q|t_1)].$$

It is crucial to note that all f -divergence measures, including the Hockey-stick and Rényi divergences, adhere to these three properties [67, Proposition 7.1, Theorem 7.2]. When the distribution of the marginal/randomness source variable T_1 is the same for P and Q , we use $D(P\|Q|t_1)$ to represent $D(P|t_1\|Q|t_1)$ for simplicity.

A. Privacy Damage of In-out Attacks

In Theorem 1, we formally demonstrate that the privacy of the shuffle-then-randomize model under in-out attack deteriorates to the local level (no privacy amplification). In Theorem 5 (see Appendix D), we reveal that the privacy of subsample-randomize-shuffle model under in-out attack degrades to the intra-batch level (amplification by $|U_k| = s$ users). Similar results applies to other shuffle model variants.

Theorem 1 (Shuffle-then-randomize model under in-out attack): Given two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$, and a protocol \mathcal{P}_{s-r} in the shuffle-then-randomize model. Let io_v denote the in-out information about user v that $io_v(k) = 1$, then for any distance measure D that satisfies the data processing inequality and separability property:

$$\begin{aligned} & D(\mathcal{P}_{s-r}(X)\|\mathcal{P}_{s-r}(X')|io_v(k) = 1) \\ & \geq \min_{z_{[0:k-1]} \in \mathcal{Z}_0 \times \dots \times \mathcal{Y}_{k-1}} D(\mathcal{R}_k(z_{[0:k-1]}, a)\|\mathcal{R}_k(z_{[0:k-1]}, b)). \end{aligned}$$

B. Privacy Damage of Message-length Attacks

Using divide-randomize-shuffle model as an example, we illustrate the destructive power of message-length attacks in Theorem 2. The victim user can only hide among users with the same message length in the same division. For the shuffle-then-randomize and subsample-randomize-shuffle models, the privacy amplification population under message-length attacks is limited to $U_{*,l}$, which represents users with identical message lengths across all rounds. When joint in-out and message-length attacks occur, the privacy amplification population further deteriorates to the same level as divide-randomize-shuffle model: $U_{k,l}$ (see Appendix F).

Theorem 2 (Divide-randomize-shuffle model under message-length attack): Consider a protocol \mathcal{P}_{d-r-s} and non-overlapping complete user divisions $\{U_k\}_{k \in [K]}$ in the divide-randomize-shuffle model, and two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$ that differ at the v -th user data. Assuming $v \in U_k$ and $len(\mathcal{R}_{(k)}(z_{(0:k-1)}, a)) \stackrel{d}{=} len(\mathcal{R}_{(k)}(z_{(0:k-1)}, b))$, let l denote the observed message-length information about user i . Define $U_{k,l}$ as the set of users with the same message length (i.e., $U_{k,l} = \{i \mid \text{for } i \in U_k \text{ and } len(\mathcal{R}_{(k)}(x_i)) = l\}$), and let $S = \{X(i)\}_{i \in U_{k,l}}$, $S' = \{X'(i)\}_{i \in U_{k,l}} \in \mathbb{X}^{|U_{k,l}|}$ denote neighboring datasets w.r.t. $U_{k,l}$ then for any distance measure D that satisfies the data processing inequality and the separability property:

$$\begin{aligned} & D(\mathcal{P}_{d-r-s}(X)\|\mathcal{P}_{d-r-s}(X')|U_{k,l}, len_v = l) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S)\|\mathcal{P}_{s-r}(S')|z_0), \end{aligned}$$

where $\mathcal{R}_{[U_{k,l}]} = \mathcal{R}_{(k)}$ are local randomizers of a shuffle-then-randomize protocol \mathcal{P}_{s-r} and z_0 is the global information.

When the condition that message lengths are distributionally equal given whether $x_v = a$ or $x_v = b$:

$$len(\mathcal{R}_{(k)}(z_{(0:k-1)}, a)) \stackrel{d}{=} len(\mathcal{R}_{(k)}(z_{(0:k-1)}, b))$$

does not hold in the theorem, there is also local privacy loss due to message length (see Section IV-E).

VI. DEFENDING AGAINST SIDE-CHANNEL ATTACKS

In this section, we present countermeasures for defending against side-channel attacks. These countermeasures give rise to new principles in the shuffle model that offer robustness to potential attacks and yield stronger privacy amplification effects. Prior to discussing the specifics, we introduce a novel model: the multinomial-randomize-shuffle (MRS) model, which possesses several advantages in decentralized settings and serves as a foundation for defense against side-channel attacks. For other variants of the shuffle model, the defense techniques are essentially the same.

A. Multinomial-randomize-shuffle Model

In this model, each user randomly selects one query k_i from $[K]$ according to a distribution $P_K \in \Delta_K$ and then responds with $\mathcal{R}_{k_i}(z_{(0:k_i-1)}, x_i)$ in the k_i -th round. The shuffler uniformly permutes messages received in the k_i -th round and

releases them to the analyzer. The query selection of each user is private and is not exposed other parties. Consequently, this model has similar privacy amplification effects as the ideal model (the privacy amplification population is n , see Theorem 4, proved in Appendix B).

Definition 11 (Multinomial-randomize-shuffle model): Let P_K denote a public distribution over $[K]$, and let k_1, \dots, k_n denote n independent samples following P_K . Define $U_k = \{i \mid \text{for } i \in [n] \text{ and } k_i = k\}$ as the k -subgroup of users. Let $\mathcal{R}_{(k)} : \mathbb{Z}_0 \times \dots \times \mathbb{Z}_{k-1} \times \mathbb{X} \rightarrow \mathbb{Y}_{(k)}$ denote the randomizer in the k -th round, where $\mathbb{Z}_k = \mathbb{Y}_{(k)}^*$ and $\mathbb{Y}_{(k)}$ is the range space of $\mathcal{R}_{(k)}$. The $z_{(0)} \in \mathbb{Z}_0$ is global information. A protocol $\mathcal{P}_{m-r-s} : \mathbb{Z}_0 \times \mathbb{X}^n \rightarrow \mathbb{Z}_0 \times \dots \times \mathbb{Z}_K$ in the MRS model proceeds as follows: given a dataset $x_{[1:n]} \in \mathbb{X}^n$, it samples $k_1, \dots, k_n \sim P_K$ to obtain $\{U_k\}_{k \in [K]}$, then sequentially compute $z_{(i)} = \mathcal{S}(\{\mathcal{R}_{(k)}(z_{(0:k-1)}, x_i)\}_{i \in U_k})$, finally outputs $z_{(0)}, \dots, z_{(K)}$ and the distribution P_K .

The MRS model has some resemblance to the random check-in approach [16]. However, the random check-in involves complex dummy and uniform selection operations on received messages in each round, and only provides privacy amplified by K users, rendering it significantly weaker than the privacy guarantees of MRS. In benign settings, MRS is also preferable to the divide-randomize-shuffle model as it offers greater potential for privacy amplification, yielding asymptotic savings of $(1 - \sqrt{1/K}) \cdot 100\%$ in budgets when all local randomizers comply with LDP. The MRS exhibits similar asymptotic privacy consumption behavior as the subsample-randomize-shuffle model [24], [42]. Assuming there are $K \approx \frac{n}{s}$ queries (i.e., one epoch in federated learning), both models consume $(\tilde{O}(\sqrt{e^{\epsilon_0}/n}), \delta)$ -DP. Since the MRS ensures each user participates in exactly one query during one epoch while the subsample-randomize-shuffle model has higher sampling variance, the constant factor in privacy loss of MRS is slightly smaller. Specifically, under the same settings as [24], [42]: $n = 60000$, $s = 1000$, $\epsilon_0 = 2$, and $K = n/s = 60$, the subsample-randomize-shuffle model consumes privacy of $(0.0367, 10^{-5})$ -DP (using the near-optimal shuffle amplification [41] and tight numerical composition with subsampling [68]), while the MRS consumes only $(0.0357, 10^{-5})$ -DP (using [41] as well). More results are listed in Table III, the MRS saves 2%-50% budget than subsampling. Moreover, the subsample-randomize-shuffle relies on additional trust in the shuffler for subsampling, whereas the MRS grants control to the user.

TABLE III
AMPLIFIED PRIVACY LEVEL COMPARISON OF MRS AND
SUBSAMPLE-RANDOMIZE-SHUFFLE MODEL ($n = 60000$, $\epsilon_0 = 2$, AND
VARYING NUMBER OF EPOCHS $E = Ks/n$).

	$\delta = 10^{-5}$			
	$E = 1$	$E = 10$	$E = 100$	$E = 500$
MRS	0.0357	0.126	0.444	1.070
s=10 ²	0.0393	0.143	0.566	1.584
s=10 ³	0.0367	0.129	0.458	1.132
	$\delta = 10^{-8}$			
	$E = 1$	$E = 10$	$E = 100$	$E = 500$
MRS	0.0571	0.190	0.635	1.485
s=10 ²	0.0629	0.211	0.769	2.037
s=10 ³	0.0593	0.194	0.651	1.551

B. Defend Against Message-length Attack

In this subsection, we present a countermeasure to defend against potential message-length attacks. Recall that message-length information possessed by adversaries threatens both local privacy and privacy amplification. To avoid local privacy loss, it is necessary to ensure that the message-length distribution is independent of the true value x_i ; to avoid degradation in privacy amplification as shown in Equation 1, it is necessary to ensure that the message-length is always the same across all users. The latter requirement is stricter as it demands identical concrete message-lengths for all users with various randomizers and true values.

Padding messages. We employ a straightforward approach, message padding (e.g., as [69] for anonymous channels), to defend against message-length attacks. Specifically, we let len_{max} denote the maximum possible message length outputted from the original local randomizers across all users:

$$len_{max} = \max_{i \in [n]} \sup_{x \in \mathbb{X}} len(\mathcal{R}_i(x)).$$

Then, for all messages released from each user, the message payload is padded to len_{max} bits.

After implementing the padding, the message length becomes global information shared among all parties. Consequently, possessing the victim's message length information no longer provides an advantage for privacy attacks.

C. Defend Against In-out Attacks

In this subsection, we present several strategies for defending against in-out attacks. For non-adaptive queries in the single-message shuffle model, there is a simple and efficient strategy: *parallelizing queries*; for adaptive queries in the single-message shuffle model, users could sacrifice communication overheads to regain privacy amplification from in-out attacks by sending *dummy messages* (a common approach in the anonymous channel literature [40]).

1) *Non-adaptive queries in the single-message shuffle model:* Many data analysis tasks involve multiple non-adaptive estimation queries. In the local model of DP, a common practice for achieving better utility (compared to dividing the privacy budget ϵ_0) is to separate the entire user population into multiple non-overlapping subsets and assign each subset to accomplish one query with the full budget ϵ_0 . For example, this approach is used in heavy hitter estimation [70], marginal queries [65], and machine learning [71]. This complies with the parallel composition theorem of differential privacy in the central model [2]. In the shuffle model, an (almost) equivalent approach is to have each user randomly choose one query among all K queries with a fixed probability distribution $P_K \in \Delta_K$, and contribute to the chosen query with the full budget [1], [43]. We illustrate this approach in Algorithm 3. Since all mechanisms M_k ($k \in [K]$) are ϵ_0 -LDP, the overall algorithm is ϵ_0 -LDP.

Given that every user follows the same distribution P_K , this implies all users are adopting an identical randomization algorithm, which ensures that privacy amplification via shuf-

Algorithm 3: Parallel local randomizer [1], [43]

Params: A distribution $P_K : [K] \mapsto [0, 1]$, base randomizers $\{\mathcal{M}_k : \mathbb{X} \mapsto \mathbb{Y}_k\}_{k \in [K]}$ each satisfies ϵ_0 -LDP.
Input: An input $x \in \mathbb{X}$.
Output: An output that satisfies ϵ_0 -LDP.

```

1 sample  $k \sim P_K$ 
2  $y \leftarrow \mathcal{M}_k(x)$ 
3 return  $y$ 

```

Algorithm 4: Rectified parallel local randomizer

Params: A distribution $P_K : [K] \mapsto [0, 1]$, base randomizers $\{\mathcal{M}_k : \mathbb{X} \mapsto \mathbb{Y}_k\}_{k \in [K]}$ each satisfies ϵ_0 -LDP, the maximum possible length len_{max} .
Input: An input $x \in \mathbb{X}$.
Output: An output that satisfies ϵ_0 -LDP.

```

1 sample  $k \sim P_K$ 
2  $y \leftarrow \mathcal{M}_k(x)$ 
3 add padding bits to  $y$  to form a  $len_{max}$ -length bit vector  $y'$ 
4 encrypts  $y'$  with the public key of the analyzer and get  $Enc_a(y')$ 
5 return  $Enc_a(y')$ 

```

fling still holds. Denote Algorithm 3 as \mathcal{R} , one straightforward conclusion is that:

$$D(S \circ \mathcal{R}(X) \parallel S \circ \mathcal{R}(X')) \leq D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X')),$$

where the shuffle-then-randomize model employs \mathcal{R} as the local randomizers.

Parallel as a cost-free defense to in-out attacks. We emphasize that parallel local randomizers are naturally immune to in-out attacks in the shuffle model. As all K queries are performed together in one single round (with the parallel randomizer in Algorithm 3), the in-out information $i_{ov} = \{1\}$ is trivial, thus maintaining the population size for shuffle privacy amplification at n . We note that parallelizing queries does not require the shuffler/server to process all messages from users within a specific short period of time but rather serves to obfuscate which specific query the victim participated in. The message sending and receiving time window can span a wide range, and the shuffler/server could process these messages in the virtual timestamp that corresponds to this time window.

Defend joint in-out & message-length attacks. We note that message-length information must be protected in the parallel local randomizer, as each base randomizer \mathcal{M}_k often has a different output space \mathbb{Y}_k , and the message length leaks information about each query the victim participated in. We present an implementation that resists joint in-out and message-length attacks in Algorithm 4, where messages are padded. Now for the privacy properties of Algorithm 4, since it can be considered as post-processing upon Algorithm 3 and both in-out & message-length information are trivial, it enjoys the same local and shuffle privacy guarantees as the shuffle-then-randomize model (with identical local randomizers).

2) *Adaptive queries in single-message shuffle model:* For adaptive queries, parallelizing queries is not applicable since the k -th query relies on previous querying results $z_{(0)}, \dots, z_{(k-1)}$. We propose contributing dummy messages at every round based on the MRS model (see Algorithm 5), at

Algorithm 5: Rectified MRS model

Params: A probability distribution $P_K : [K] \mapsto [0, 1]$, adaptive local randomizers $\{\mathcal{R}_{(k)}\}_{k \in [K]}$, the maximum possible length len_{max} , global information $z_{(0)}$.
Input: Inputs $x_1, \dots, x_n \in \mathcal{X}$ from n users.
Output: The querying results of K adaptive queries.

```

1  $\triangleright$  Sample participation choices on the user side
2 for users  $i \in [n]$  do
3    $\lfloor$  sample  $k_i \sim P_K$ 
4  $\triangleright$  Run randomization & shuffling
5 for  $k \in [K]$  do
6   for users  $i \in [n]$  do
7      $\triangleright$  Randomize, pad, and encrypt on the user side
8     if  $k = k_i$  then
9        $y_{i,k} \leftarrow \mathcal{R}_{(k)}(z_{(0:k-1)}, x_i)$ 
10    else
11       $\lfloor$  let  $y_{i,k}$  be an empty message
12      pad  $y_{i,k}$  to form a  $len_{max}$ -length bit vector  $y'_{i,k}$ 
13      encrypts  $y'_{i,k}$  with the analyzer's public key to get  $Enc_a(y'_{i,k})$ 
14  $\triangleright$  Uniform-randomly permute on the shuffler side
15  $E_k = \mathcal{S}(Enc_a(y'_{1,k}), \dots, Enc_a(y'_{n,k}))$ 
16  $\triangleright$  Decrypt and analyze on the server side
17  $Y'_k = \{Dec_a(ciphertext) \mid ciphertext \in E_k\}$ 
18  $Y_k = \{y \mid y \in Y'_k \text{ and } y \text{ is not empty}\}$ 
19  $z_{(k)} = \mathcal{A}_k(Y_k)$ 
20 return  $z_{(0)}, \dots, z_{(K)}$ 

```

the cost of $(K-1)$ -times more communication overhead. We then show how to trade-off communication and amplification.

Sending dummy messages. To obscure the information about which round a user participated in and prevent adversaries from obtaining meaningful in-out information, we let every user contribute message(s) at all K rounds. In the true participating round k_i (sampled as in the normal MRS model), user i contributes a true message $\mathcal{R}_{(k_i)}(z_{(0:k-1)}, x_i)$ (line 9); in the other $K-1$ rounds, user i contributes a dummy/empty message (line 11). To further prevent adversaries from launching message-length attacks, all messages are padded (line 12) and encrypted before transmission.

Since every user participates in all rounds, the in-out information $i_{ov} = [1, \dots, 1]$ becomes trivial in the new model presented in Algorithm 5; since all messages have the same length, there is no degradation due to side-channel message-length information. Moreover, the observed shuffled messages E_1, \dots, E_K can be viewed as adding padded empty messages at every round to the outputs of \mathcal{P}_{m-r-s} according to the total number of users n (a public piece of information). Therefore, the new model resists in-out & message-length attacks and has the same privacy amplification effects as the normal MRS model in Theorem 4.

A practical issue in the new model is that the total messages grows to K times when compared to the original model. This might pose significant communication burdens on users with scarce resources (e.g., in mobile devices). To balance communication costs and privacy amplification, we propose a flexible, generalized participation model extending the multinomial approach.

Algorithm 6: Rectified bin-randomize-shuffle model

Params: Participation bins $Q_{(1:M)}$, probability distributions $P_{(m)} : Q_{(m)} \mapsto [0, 1]$ for $m \in [M]$, adaptive local randomizers $\{\mathcal{R}_{(k)}\}_{k \in [K]}$, the maximum possible length len_{max} , global information $z_{(0)}$.

Input: Inputs $x_1, \dots, x_n \in \mathcal{X}$ from n users.

Output: The querying results of K adaptive queries.

```

1 ▷ Compute participation choices on the user side
2 for users  $i \in [n]$  do
3   choose  $m_i \in [M]$  with an arbitrary (personalized) rule
4    $k_i \sim P_{(m_i)}$ 
5 ▷ Run randomization & shuffling
6 for  $k \in [K]$  do
7   for users  $i \in [n]$  do
8     ▷ Randomize, pad, and encrypt on the user side
9     if  $k \in Q_{(m_i)}$  then
10      if  $k = k_i$  then
11         $y_{i,k} \leftarrow \mathcal{R}_{(k)}(z_{(0:k-1)}, x_i)$ 
12      else
13        let  $y_{i,k}$  be an empty message
14      pad  $y_{i,k}$  to form a  $len_{max}$ -length bit vector  $y'_{i,k}$ 
15      encrypts  $y'_{i,k}$  with the analyzer's public key to
      get  $Enc_a(y_{i,k})$ 
16    ▷ Uniform-randomly permute on the shuffler side
17     $E_k = \mathcal{S}(\{Enc_a(y'_{i,k})\}_{i \in U_{i(k)}})$ 
18    ▷ Decrypt and analyze on the server side
19     $Y'_k = \{Dec_a(ciphertext) \mid ciphertext \in E_k\}$ 
20     $Y_k = \{y \mid y \in Y'_k \text{ and } y \text{ is not empty}\}$ 
21     $z_{(k)} = \mathcal{A}_k(Y_k)$ 
22 return  $z_{(0)}, \dots, z_{(K)}$ 

```

Bin participation paradigm. The model first defines a series of bar points $\{b_0, b_1, \dots, b_M\}$ in $[K + 1]$. Specifically, $b_0 \equiv 1$, $b_M \equiv K + 1$, and $b_m < b_{m+1}$ for $m \in [M - 1]$. The model then defines a set $Q = \{Q_{(1)}, \dots, Q_{(M)}\}$ of non-overlapping, consecutive, and complete subsets (bins) of $[K]$ where $Q_{(m)} = [b_{m-1}, b_m - 1]$ for $m \in [M]$. It is obvious that $Q_{(m)} \neq \emptyset$, $Q_{(m)} \cap Q_{(m')} = \emptyset$ holds for all $m, m' \in [M]$ when $m \neq m'$, and $Q_{(1)} \cup Q_{(2)} \cup \dots \cup Q_{(M)} = [K]$. For each $Q_{(m)} \in Q$, it is associated with a probability distribution $P_{(m)} : Q_{(m)} \mapsto [0, 1] \in \Delta_{|Q_{(m)}|}$. In the binned participating model, each user $i \in [n]$ select $m_i \in [M]$ with an arbitrary rule. For example, if a user becomes available online after round k , the user might select a bin $Q_{(m_i)}$ such that $Q_{(m_i)} \subseteq [k + 1 : K]$; the user might also arbitrary-randomly select a bin $Q_{(m_i)}$ from Q . After choosing the bin $Q_{(m_i)}$, the user i samples a true participation round k_i from $Q_{(m_i)}$ according to $P_{(m_i)}$. The corresponding shuffle model with the binned participation paradigm is termed the bin-randomize-shuffle model. Specifically, when $Q = [n]$, the bin-randomize-shuffle model is equivalent to the MRS model.

We let $U_{(m)} \subseteq [n]$ denote the users that selected bin $Q_{(m)}$ in the bin-randomize-shuffle model, $U_k \subseteq [n]$ denote the users that selected query k , and let $l_{(k)}$ denote the bin where k belongs to (i.e., $k \in Q_{(l_{(k)})}$). Then, by combining the defensive techniques in the MRS model, we show that user $i \in U_{(m)}$ can hide among $U_{(m)}$, even under in-out and message-length attacks. We present the overall procedure in Algorithm 6, and

formally state the privacy amplification guarantee in Theorem 3 (see Appendix G for proof). Since $U_{(m)}$ lies between U_k and $[n]$, the privacy amplification and communication costs can be flexibly traded off.

Theorem 3 (Rectified bin-randomize-shuffle model): Given a protocol \mathcal{P}_{b-r-s} , non-overlapping complete bins $Q = \{Q_{(1)}, \dots, Q_{(M)}\}$, and non-overlapping complete user divisions $\{U_{(m)}\}_{m \in [M]}$ in the binned-randomize-shuffle model (as in Algorithm 6), and two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$. Let l denote the message-length information about user v that $v \in U_k \subseteq U_{(m)}$, and let $S = \{X(i')\}_{i' \in U_{(m)}}$, $S' = \{X'(i')\}_{i' \in U_{(m)}} \in \mathbb{X}^{|U_{(m)}|}$ denote neighboring datasets w.r.t. $U_{(m)}$ then for any distance measure D that satisfies the data processing inequality and separability property:

$$D(\mathcal{P}_{b-r-s}(X) \parallel \mathcal{P}_{b-r-s}(X') \mid U_{(m)}, len_v = l) \leq \max_{z_0} D(\mathcal{P}_{m-r-s}(S) \parallel \mathcal{P}_{m-r-s}(S') \mid z_0),$$

where $\{\mathcal{R}_{(k')}\}_{k' \in Q_{m_i}}$ are the local randomizers of a multinomial-randomize-shuffle protocol \mathcal{P}_{m-r-s} that has $|Q_m|$ rounds, z_0 is the global information in \mathcal{P}_{m-r-s} and $P_{(m)}$ is the query selection distribution of the m -th bin.

D. Defend Against Message-cardinality Attacks

Section IV-C highlights that message-cardinality attacks may result in significant privacy loss for some SOTA protocols [17], [19], [46], [47]. In this section, we demonstrate the difficulty and feasibility in securing the Δ -summation protocol [19] against message-cardinality attacks.

Recall that in the Δ -summation protocol (see Algorithms 1 and 2) with $\Delta = 1$, if a user has $x_i = 1$, they send a 1 and a random number of $-1, 1$ messages to the shuffler. Conversely, if a user has $x_i = 0$, they send only a random number of $-1, 1$ messages following the same distribution as when $x_i = 1$. A straightforward remedy would involve sending an extra 0 (or any other stub/dummy message) to the shuffler when $x_i = 0$ to ensure identical message cardinality distributions for both $x_i = 1$ and $x_i = 0$. However, while this addresses the message-cardinality issue, the number of 0s or stubs observed by the analyzer (i.e., potential adversaries) in shuffled messages directly reveals the number of users with 0, leading to an infinite differential privacy loss.

One feasible approach is transforming the Δ -summation problem into a $(\Delta + 1)$ -summation problem while utilizing the original protocol as a base. In this approach, given user data $x_i \in [\Delta]$, we first transform it to $x_i + 1 \in [\Delta + 1]$. These transformed values are then fed into a $(\Delta + 1)$ -summation protocol from [19]. After retrieving the noisy summation s from the protocol, we obtain the final unbiased estimator of $\sum_{i \in [n]} x_i$ by subtracting n from s . By incorporating the transformation step where we only input $x_i + 1$, which is always greater than 0, the local randomizer in Algorithm 1 consistently sends $x_i + 1$, effectively eliminating the message-cardinality issue. Since the user population size n is a publicly known parameter and each user data point increases by 1, we can still recover the unbiased noisy summation. The only additional cost incurred is the increased message complexity

associated with a $(\Delta+1)$ -summation protocol when compared to a Δ -summation protocol.

It is important to note that this defense approach cannot be generalized to other vulnerable protocols in [17], [46], [47] that only take binary inputs.

VII. DISCUSSIONS

In this section, we summary new privacy requirements and novel principles for the shuffle model, and discuss their implications to more broader DP-Cryptography systems.

A. Curious Shuffle Model

Keep in mind that privacy attackers can easily carry out side-channel attacks on the shuffle model. A typical example of such attackers is an honest-but-*curious* shuffler. If we let $\text{SideInfo}(i)$ represent the information that privacy attackers have about a user $i \in [n]$, we can then define a side-channel-resistant version of DP in the shuffle model like this:

Definition 12 (DP in the curious shuffle model): A protocol $\mathcal{P} = (\{\mathcal{R}_i\}_{i \in [n]}, \mathcal{A})$ satisfies (ϵ, δ) -differential privacy in the curious shuffle model iff for all neighboring datasets X and $X' \in \mathbb{X}^n$, the $(\mathcal{S} \circ \mathcal{R}_{[n]}(X), \text{SideInfo}(i)_{i \in [n]})$ and $(\mathcal{S} \circ \mathcal{R}_{[n]}(X'), \text{SideInfo}(i)_{i \in [n]})$ are (ϵ, δ) -indistinguishable.

Usually, $\text{SideInfo}(i)$ includes $io(i)$, $num(\mathcal{R}_i)$ and $len(\mathcal{R}_i)$ as defined in Section IV-B. In this case, our rectified proposals in Section VI maintain privacy amplification in the curious shuffle model.

B. New Principles in the Shuffle Model

To maximize privacy amplification effects and minimize side-channel risks, we summarize new principles for the (curious) shuffle model:

(a) Pad every message. As demonstrated in Section IV-E, message-length information can significantly compromise privacy. A simple but effective solution is to pad each message to at least the maximum possible message length [69], len_{max} . Consequently, the message length becomes a public information in the system, thus avoid privacy degradation even when adversaries possess message length information.

(b) Parallelize queries. When handling multiple queries with no sequential dependence, the best practice in the shuffle model is to pack them into one parallel query (see Section VI-C1). This allows users to benefit from privacy amplification over the full population without sending dummy messages.

(c) Contribute dummy messages. As for adaptive queries, existing shuffle models fail to fully amplify privacy across rounds when there are side-channel attacks. An effective mitigation approach is to have each user contribute extra dummy messages [40] to the shuffler at every round (as shown in Algorithm 5), so that users can hide among the overall population. Alternatively, users may adopt the bin-randomize-shuffle model (see Algorithm 6) to flexibly control the trade-off between communication overheads, privacy amplification effects, and utility.

Combining these principles enables users to achieve optimized privacy-utility-communication trade-offs.

C. New Challenges in DP-Crypto Systems

Stemming from this study, while countermeasures from the cryptography literature (e.g., padding bits and sending dummy messages) offer valuable insights, they are not universally applicable to defending against side-channel attacks in the shuffle DP model (see Section VI-D) and may be inefficient (see Section VI-C), due to differences in functionalities and privacy goals. Our research pivots to a crucial frontier: the synthesis of information-theoretic privacy tools and cryptographic methods (e.g., anonymous channels, oblivious shuffling, order preserving encryption, and approximate homomorphic encryption). This merger seeks not only to balance privacy and utility/efficiency but also to illuminate and address novel attack vectors. These attacks, within the DP context, beckon comprehensive analyses and the genesis of innovative defenses.

VIII. CONCLUSION

This study presents a novel identification of communication side-channel attacks associated with the shuffle model of differential privacy. We classify these attacks into three categories: in-out, message-length, and message-cardinality attacks. Additionally, we empirically and theoretically investigate the resulting degradation of privacy due to these attacks. Our findings reveal that these attacks cause a significant or even infinity increase in privacy loss, potentially nullifying the benefits of privacy amplification through shuffling. To counteract these vulnerabilities, we propose two new variants of the shuffle model: the MRS model and the bin-randomize-shuffle model. We introduce new principles within these models, such as parallelizing queries, padding messages, and sending dummy messages. As a result, the privacy amplification effects are redeemed with minimal additional costs. Furthermore, these new models are applicable not only for defending against attacks but also in benign environments, resulting in significantly stronger privacy amplification effects compared to existing models.

REFERENCES

- [1] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, "Amplification by shuffling: From local to central differential privacy via anonymity," *SODA*, 2019.
- [2] C. Dwork, "Differential privacy," in *ICALP*. Springer, 2006.
- [3] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Local privacy and statistical minimax rates," in *FOCS*. IEEE, 2013.
- [4] S. Wang, L. Huang, Y. Nie, X. Zhang, P. Wang, H. Xu, and W. Yang, "Local differential private data aggregation for discrete distribution estimation," *IEEE Transactions on Parallel and Distributed Systems*, vol. 30, no. 9, pp. 2046–2059, 2019.
- [5] S. Wang, L. Huang, Y. Nie, P. Wang, H. Xu, and W. Yang, "Privset: Set-valued data analyses with locale differential privacy," in *INFOCOM*. IEEE, 2018.
- [6] S. Wang, Y. Qian, J. Du, W. Yang, L. Huang, and H. Xu, "Set-valued data publication with local privacy: tight error bounds and efficient mechanisms," *VLDB*, 2020.
- [7] S. Wang, Y. Peng, K. Chen, and W. Yang, "Optimal locally private data stream analytics," in *INFOCOM*. IEEE, 2024.
- [8] S. Wang, X. Luo, Y. Qian, J. Du, W. Lin, and W. Yang, "Analyzing preference data with local privacy: Optimal utility and enhanced robustness," *IEEE Transactions on Knowledge and Data Engineering*, 2022.

- [9] S. Wang, Y. Li, Y. Zhong, K. Chen, X. Wang, Z. Zhou, F. Peng, Y. Qian, J. Du, and W. Yang, "Locally private set-valued data analyses: Distribution and heavy hitters estimation," *IEEE Transactions on Mobile Computing*, 2023.
- [10] R. Dingleline, N. Mathewson, and P. Syverson, "Tor: The second-generation onion router," Naval Research Lab Washington DC, Tech. Rep., 2004.
- [11] V. Costan and S. Devadas, "Intel sgx explained," *Cryptology ePrint Archive*, 2016.
- [12] A. Bittau, Ú. Erlingsson, P. Maniatis, I. Mironov, A. Raghunathan, D. Lie, M. Rudominer, U. Kode, J. Tinnes, and B. Seefeld, "Prochlo: Strong privacy for analytics in the crowd," in *ACM SOSP*, 2017.
- [13] I. Abraham, B. Pinkas, and A. Yanai, "Blinder—scalable, robust anonymous committed broadcast," in *CCS*. ACM, 2020.
- [14] V. Feldman, A. McMillan, and K. Talwar, "Stronger privacy amplification by shuffling for rényi and approximate differential privacy," in *SODA*. SIAM, 2023.
- [15] B. Balle, J. Bell, A. Gascón, and K. Nissim, "The privacy blanket of the shuffle model," *CRYPTO*, 2019.
- [16] B. Balle, P. Kairouz, B. McMahan, O. D. Thakkar, and A. Thakurta, "Privacy amplification via random check-ins," *NeurIPS*, 2020.
- [17] B. Ghazi, R. Kumar, P. Manurangsi, and R. Pagh, "Private counting from anonymous messages: Near-optimal accuracy with vanishing communication overhead," in *ICML*. PMLR, 2020.
- [18] V. Balcer, A. Cheu, M. Joseph, and J. Mao, "Connecting robust shuffle privacy and pan-privacy," in *SODA*. SIAM, 2021.
- [19] B. Ghazi, R. Kumar, P. Manurangsi, R. Pagh, and A. Sinha, "Differentially private aggregation in the shuffle model: Almost central accuracy in almost a single message," in *ICML*. PMLR, 2021.
- [20] S. Wang, J. Li, Y. Qian, J. Du, W. Lin, and W. Yang, "Hiding numerical vectors in local private and shuffled messages," in *IJCAI*, 2021.
- [21] A. Cheu and M. Zhilyaev, "Differentially private histograms in the shuffle model from fake users," in *IEEE S&P*, 2022.
- [22] S. Wang, X. Luo, Y. Qian, Y. Zhu, K. Chen, Q. Chen, B. Xin, and W. Yang, "Shuffle differential private data aggregation for random population," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 5, pp. 1667–1681, 2023.
- [23] S. Wang, S. Yu, X. Ren, J. Li, Y. Li, W. Yang, and H. Yan, "Differentially private numerical vector analyses in the local and shuffle model," *IEEE Transactions on Dependable and Secure Computing*, 2023.
- [24] A. Girgis, D. Data, S. Diggavi, P. Kairouz, and A. T. Suresh, "Shuffled model of differential privacy in federated learning," in *AISTATS*, 2021.
- [25] J. Tenenbaum, H. Kaplan, Y. Mansour, and U. Stemmer, "Differentially private multi-armed bandits in the shuffle model," *NeurIPS*, 2021.
- [26] A. Lowy and M. Razaviyayn, "Private federated learning without a trusted server: Optimal algorithms for convex losses," *arXiv preprint arXiv:2106.09779*, 2021.
- [27] A. Cheu, M. Joseph, J. Mao, and B. Peng, "Shuffle private stochastic convex optimization," in *ICLR*, 2022.
- [28] S. R. Chowdhury and X. Zhou, "Shuffle private linear contextual bandits," in *ICML*. PMLR, 2022.
- [29] E. Garcelon, K. Chaudhuri, V. Perchet, and M. Pirotta, "Privacy amplification via shuffling for linear contextual bandits," in *ALT*, 2022.
- [30] F. Li, X. Zhou, and B. Ji, "Differentially private linear bandits with partial distributed feedback," in *2022 20th International Symposium on Modeling and Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt)*. IEEE, 2022, pp. 41–48.
- [31] Apple and Google, "Exposure notification privacy-preserving analytics white paper," 2021, https://covid19-static.cdn-apple.com/applications/covid19/current/static/contact-tracing/pdf/ENPA_White_Paper.pdf.
- [32] Y.-S. Shiu, S. Y. Chang, H.-C. Wu, S. C.-H. Huang, and H.-H. Chen, "Physical layer security in wireless networks: A tutorial," *IEEE wireless Communications*, vol. 18, no. 2, pp. 66–74, 2011.
- [33] Y. Wan, K. Xu, F. Wang, and G. Xue, "Characterizing and mining traffic patterns of iot devices in edge networks," *IEEE Transactions on Network Science and Engineering*, vol. 8, no. 1, pp. 89–101, 2020.
- [34] S. Feghhi and D. J. Leith, "A web traffic analysis attack using only timing information," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 8, pp. 1747–1759, 2016.
- [35] J. Deng, R. Han, and S. Mishra, "Countermeasures against traffic analysis attacks in wireless sensor networks," in *SECURECOMM*. IEEE, 2005.
- [36] F. Zhang, W. He, and X. Liu, "Defending against traffic analysis in wireless networks through traffic reshaping," in *ICDCS*. IEEE, 2011.
- [37] S. J. Murdoch and G. Danezis, "Low-cost traffic analysis of tor," in *IEEE S&P*, 2005.
- [38] L. Overlier and P. Syverson, "Locating hidden servers," in *IEEE S&P*, 2006.
- [39] G. Danezis and A. Serjantov, "Statistical disclosure or intersection attacks on anonymity systems," in *Information Hiding: 6th International Workshop*. Springer, 2005, pp. 293–308.
- [40] O. Berthold and H. Langos, "Dummy traffic against long term intersection attacks," in *Privacy Enhancing Technologies: Second International Workshop, PET 2002 San Francisco, CA, USA, April 14–15, 2002 Revised Papers 2*. Springer, 2003, pp. 110–128.
- [41] V. Feldman, A. McMillan, and K. Talwar, "Hiding among the clones: A simple and nearly optimal analysis of privacy amplification by shuffling," in *IEEE FOCS*, 2022.
- [42] A. Girgis, D. Data, and S. Diggavi, "Rényi differential privacy of the subsampled shuffle model in distributed learning," *NeurIPS*, 2021.
- [43] S. Wang, Y. Peng, J. Li, Z. Wen, Z. Li, S. Yu, D. Wang, and W. Yang, "Privacy amplification via shuffling: Unified, simplified, and tightened," *Proceedings of the VLDB Endowment*, vol. 17, no. 8, pp. 1870–1883, 2024.
- [44] Z. Weinberg, J. Wang, V. Yegneswaran, L. Briesemeister, S. Cheung, F. Wang, and D. Boneh, "Stegotorus: a camouflage proxy for the tor anonymity system," in *CCS*, 2012.
- [45] D. L. Chaum, "Untraceable electronic mail, return addresses, and digital pseudonyms," *Communications of the ACM*, vol. 24, no. 2, pp. 84–90, 1981.
- [46] B. Ghazi, N. Golowich, R. Kumar, P. Manurangsi, R. Pagh, and A. Velingker, "Pure differentially private summation from anonymous messages," in *1st Conference on Information-Theoretic Cryptography (ITC 2020)*, 2020.
- [47] V. Balcer and A. Cheu, "Separating local & shuffled differential privacy via histograms," in *1st Conference on Information-Theoretic Cryptography*, 2020.
- [48] J. M. Abowd, "The us census bureau adopts differential privacy," in *ACM SIGKDD*, 2018.
- [49] Ú. Erlingsson, V. Pihur, and A. Korolova, "Rappor: Randomized aggregatable privacy-preserving ordinal response," in *CCS*. ACM, 2014.
- [50] I. Mironov, "On significance of the least significant bits for differential privacy," in *ACM CCS*, 2012.
- [51] J. Jin, E. McMurtry, B. I. Rubinfeld, and O. Ohrimenko, "Are we there yet? timing and floating-point attacks on differential privacy systems," in *SP*. IEEE, 2022.
- [52] A. Haeberler, B. C. Pierce, and A. Narayan, "Differential privacy under fire," in *20th USENIX Security Symposium (USENIX Security 11)*, 2011.
- [53] A. Cheu, A. Smith, and J. Ullman, "Manipulation attacks in local differential privacy," in *SP*. IEEE, 2021.
- [54] T. Wang, B. Ding, M. Xu, Z. Huang, C. Hong, J. Zhou, N. Li, and S. Jha, "Improving utility and security of the shuffler-based differential privacy," *Proceedings of the VLDB Endowment*, vol. 13, no. 13, pp. 3545–3558, 2020.
- [55] A. Cheu, A. Smith, J. Ullman, D. Zeber, and M. Zhilyaev, "Distributed differential privacy via shuffling," *EUROCRYPT*, 2019.
- [56] B. Balle, G. Barthe, and M. Gaboardi, "Privacy amplification by subsampling: Tight analyses via couplings and divergences," *NeurIPS*, 2018.
- [57] Q. Wang, A. Yahyavi, B. Kemme, and W. He, "I know what you did on your smartphone: Inferring app usage over encrypted data traffic," in *2015 IEEE conference on communications and network security (CNS)*. IEEE, 2015, pp. 433–441.
- [58] J. S. Atkinson, J. E. Mitchell, M. Rio, and G. Match, "Your wifi is leaking: What do your mobile apps gossip about you?" *Future Generation Computer Systems*, vol. 80, pp. 546–557, 2018.
- [59] E. Papadogiannaki and S. Ioannidis, "A survey on encrypted network traffic analysis applications, techniques, and countermeasures," *ACM Computing Surveys (CSUR)*, vol. 54, no. 6, pp. 1–35, 2021.
- [60] D. Naboulsi, M. Fiore, S. Ribot, and R. Stanica, "Large-scale mobile traffic analysis: a survey," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 124–161, 2015.
- [61] J. Ren, A. Rao, M. Lindorfer, A. Legout, and D. Choffnes, "Recon: Revealing and controlling pii leaks in mobile network traffic," in *ACM MobiSys*, 2016.
- [62] M. Conti, Q. Q. Li, A. Maragno, and R. Spolaor, "The dark side (-channel) of mobile devices: A survey on network traffic analysis," *IEEE communications surveys & tutorials*, vol. 20, no. 4, pp. 2658–2713, 2018.
- [63] L. Chen, B. Ghazi, R. Kumar, and P. Manurangsi, "On distributed differential privacy and counting distinct elements," in *12th Innovations in Theoretical Computer Science Conference (ITCS 2021)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2021.

- [64] G. Cormode, T. Kulkarni, and D. Srivastava, "Answering range queries under local differential privacy," *Proceedings of the VLDB Endowment*, vol. 12, no. 10, pp. 1126–1138, 2019.
- [65] —, "Marginal release under local differential privacy," in *SIGMOD*. ACM, 2018.
- [66] X. Gu, M. Li, Y. Cheng, L. Xiong, and Y. Cao, "PCKV: Locally differentially private correlated key-value data collection with optimized utility," *USENIX Security*, 2020.
- [67] Y. Polyanskiy and Y. Wu, *Information Theory: From Coding to Learning*. Cambridge University Press, 2022, <https://people.lids.mit.edu/yp/homepage/data/itbook-export.pdf>.
- [68] A. Koskela, M. A. Heikkilä, and A. Honkela, "Numerical accounting in the shuffle model of differential privacy," *Transactions on Machine Learning Research*, 2023.
- [69] S. Yu, G. Zhao, W. Dou, and S. James, "Predicted packet padding for anonymous web browsing against traffic analysis attacks," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 4, pp. 1381–1393, 2012.
- [70] R. Bassily, K. Nissim, U. Stemmer, and A. Guha Thakurta, "Practical locally private heavy hitters," *NeurIPS*, 2017.
- [71] S. Fletcher and M. Z. Islam, "Decision tree classification with differential privacy: A survey," *ACM Computing Surveys (CSUR)*, vol. 52, no. 4, pp. 1–33, 2019.

Shaowei Wang received the PhD degree in computer science from the University of Science and Technology of China (USTC), in 2019. He is an associate professor in Institute of Artificial Intelligence at Guangzhou University. His research interests are data privacy, federated learning and recommendation systems. He has published over 40 papers on top-tier conferences and journals, such as NeurIPS, VLDB, INFOCOM, ICDE, IEEE TIFS, and IEEE TDSC.

Jin Li is a master candidate at Guangzhou University. His research interests are data privacy and federated learning.

Changyu Dong received the PhD degree from Imperial College London. He is currently a professor at the Institute of Artificial Intelligence, Guangzhou University. He has authored more than 50 publications in international journals and conferences. His research interests include applied cryptography, data privacy, AI security, and blockchain. His recent work focuses mostly on designing practical secure computation protocols. The application domains include secure cloud computing and privacy-preserving data mining.

Jin Li received the PhD degree in information security from Sun Yat-sen University, Guangzhou, China, in 2007. He is currently a professor and vice dean with the School of Computer Science, Guangzhou University. His research interests include design of secure protocols in cloud computing and cryptographic protocols. He has published more than 100 papers in top-tier international conferences and journals, including IEEE INFOCOM, IEEE Transactions on Information Forensics and Security and ESORICS etc. His work has been cited more than 10,000 times at Google Scholar and the H-Index is 34. He received NSFC Outstanding Youth Foundation, in 2017.

Zhili Zhou received the PhD degrees in Computer Application at the School of Information Science and Engineering from Hunan University, in 2014. He is currently a professor with Institute of Artificial Intelligence, Guangzhou University. His current research interests include Multimedia Security, Artificial Intelligence Security and Information Hiding. He has been selected as "World's Top 2% Scientists" from 2019 to 2022 by Stanford University and Elsevier. He received ACM Rising Star Award and got Guangdong Natural Science Funds for Distinguished Young Scholar.

Di Wang (Member, IEEE) received the Ph.D. degree in computer science from The State University of New York (SUNY) at Buffalo. He is currently an Assistant Professor of computer science and a Faculty Member of statistics with the Division of Computer, Electrical and Mathematical Sciences and Engineering (CEMSE), King Abdullah University of Science and Technology (KAUST). His research interests include trustworthy machine learning, machine learning theory, and AI for science.

Zikai Wen received his PhD degree in computer science from Cornell University, in 2021. He is a postdoctoral researcher at the department of computer science, Virginia Tech. His research homes on two critical challenges: (1) safeguarding usable privacy and security in AI interactions and (2) dismantling barriers to AI accessibility for learners with neurodevelopmental disabilities.

APPENDIX A

LOCAL PRIVACY LOSS DUE TO MESSAGE LENGTH

In some commonly-used local randomizers (e.g., RAPPOR [49], key-value data randomizer [66]), their outputs are often compressed (e.g., via list representation of sparse vector) before encryption for transmission efficiency as in [27]. The compressed message length might probabilistically reveal information about the message value, thus incur severe local privacy loss. For instance, in RAPPOR, the secret value x_i is hashed into a Bloom filter of length b using h hash functions $\{H_j\}_{j \in [h]}$, and each bit is randomly flipped with probability $p \in [0, 0.5]$. For transmission efficiency, the randomized Bloom filter is compressed as a list of non-zero indexes. Meanwhile, the length of the list might correlates with the input value. Considering $x, x' \in \mathbb{X}$ such that $1 = \#\{H_j(x)\}_{j \in [h]}$ and $h = \#\{H_j(x')\}_{j \in [h]}$, then their probability distributions of number of 1s (i.e., the length of the index list) in the randomized Bloom filter largely differ. This induces local privacy loss of $(h-1)\log((1-p)/p) = \frac{h-1}{2h}\epsilon_0$ (refer to Definition 13), where ϵ_0 is the local budget and is relatively large in the shuffle model.

Definition 13 (*Local privacy loss due to message length*): Considering a local randomizer $\mathcal{R} : \mathbb{X} \mapsto \mathbb{Y}$, let $\text{len}(\mathcal{R}(x))$ denote the number of bits of $\mathcal{R}(x)$ on a given input $x \in \mathbb{X}$. The local privacy loss of message length information is (ϵ, δ) if for some $x, x' \in \mathbb{X}$:

$$D_{\epsilon^\epsilon}(\text{len}(\mathcal{R}(x)) \parallel \text{len}(\mathcal{R}(x'))) \geq \delta.$$

APPENDIX B

PRIVACY AMPLIFICATION OF MULTINOMIAL-RANDOMIZE-SHUFFLE MODEL

Theorem 4 (*Privacy amplification of multinomial-randomize-shuffle model*): Given a protocol \mathcal{P}_{m-r-s} and sampled non-overlapping complete user divisions $\{U_k\}_{k \in [K]}$ in multinomial-randomize-shuffle model, and two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$ that differ at the v -th user data, then for any distance measure D that satisfies data processing inequality:

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \parallel \mathcal{P}_{m-r-s}(X') \mid \{U_k\}_{k \in [K]}) \\ & \leq D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X')) \end{aligned}$$

where $\mathcal{R}_{[n]} = \mathcal{R}_{(1)}^{|U_1|} \times \dots \times \mathcal{R}_{(K)}^{|U_K|}$ are the local randomizers of \mathcal{P}_{s-r} in the shuffle-then-randomize model, and the global information in \mathcal{P}_{s-r} is the same as in \mathcal{P}_{m-r-s} .

We utilize the data processing inequality to prove the theorem. Specifically, we show there exists a shuffle-then-randomize protocol \mathcal{P}_{s-r} (where $\mathcal{R}_{[n]} \in \{\mathcal{R}_{(1)}, \dots, \mathcal{R}_{(K)}\}^n$) and a post-processing function, such that the output from MRS model (distributionally) equals to post-processed output of \mathcal{P}_{s-r} .

Consider a special protocol \mathcal{P}_{s-r} in the shuffle-then-randomize model that $\mathcal{R}_i \equiv \mathcal{R}_{(k)}$ for $i \in [\sum_{k \in [k-1]} |U_k| : \sum_{k' \in [k]} |U_k|]$ and the global information z_0 is the same as in \mathcal{P}_{m-r-s} . Furthermore, we assume there is no adaptivity within U_k (the $k \in [K]$), meaning that every $\mathcal{R}_i = \mathcal{R}_{(k)}$ are independent from $\{z_{i'}\}_{i' \in U_k \text{ and } i' \leq i}$. Let $\mathbf{z} = \{z_0, \dots, z_n\}$ denote the output variables of such an special algorithm. Since all users follow the same multinomial distribution P_K , and

both \mathcal{P}'_{s-r} and \mathcal{P}_{m-r-s} uses uniform-random shuffling, due to the uniformity of all $\{x_1, \dots, x_n\}$ in the input, we have the output distribution of \mathbf{z} from the \mathcal{P}_{s-r} equal to the \mathcal{P}_{m-r-s} with the same input. Therefore, according to the post-processing inequality (with an post-processing function of identical map), we have the conclusion.

We examine the conditional cases of the multinomial-randomize-shuffle model in Theorem 4, wherein the sub-population sizes for all rounds $\{|U_k|\}_{k \in [K]}$ are fixed. Regarding non-conditional case $D(\mathcal{P}_{m-r-s}(X) \parallel \mathcal{P}_{m-r-s}(X'))$, one can effortlessly employ the conditioning increasing property of divergence measures to derive:

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \parallel \mathcal{P}_{m-r-s}(X')) \\ & \leq \mathbb{E}_{\{U_k\}_{k \in [K]} \sim \text{multinomial}(n, P_K)} D(\mathcal{P}_{m-r-s}(X) \parallel \mathcal{P}_{m-r-s}(X') \mid \{U_k\}_{k \in [K]}). \end{aligned}$$

APPENDIX C

PROOF OF THEOREM 1 ON SHUFFLE-THEN-RANDOMIZE MODEL UNDER IN-OUT ATTACKS

To prove that:

$$\begin{aligned} & D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X') \mid i_{O_v}(k) = 1) \\ & \geq \min_{z_{[0:k-1]} \in \mathbb{Z}_0 \times \dots \times \mathbb{Y}_{k-1}} D(\mathcal{R}_k(z_{[0:k-1]}, a) \parallel \mathcal{R}_k(z_{[0:k-1]}, b)), \end{aligned}$$

we consider the output of \mathcal{P}_{s-r} with fixed variables $i_{O_v}(k) = 1$ and utilize the data processing inequality. We define a post-processing function over the output of \mathcal{P}_{s-r} :

- (1): Remove the $z_{k'}$ for $k' \in [k+1 : K]$ from the output list;
- (2): Return $z_0 = z_{[0:k-1]}$ and z_k .

When $z_{[0:k-1]}$ is fixed, given X or X' , the output distribution of the above (2) is equivalent to $\mathcal{R}_k(z_{[0:k-1]}, a)$ or $\mathcal{R}_k(z_{[0:k-1]}, b)$, respectively. Therefore, according to the data processing inequality, we have:

$$\begin{aligned} & D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X') \mid i_{O_v}(k) = 1, z_{[0:k-1]} = z) \\ & \geq D(\mathcal{R}_k(z_{[0:k-1]}, a) \parallel \mathcal{R}_k(z_{[0:k-1]}, b)). \end{aligned}$$

Further since X and X' differ only at x_i and user i appears at the k -th round, in two independent runs, $\mathcal{P}_{s-r}(X) \mid i_{O_v}(k) = 1, z_{[0:k-1]}$ and $\mathcal{P}_{s-r}(X') \mid i_{O_v}(k) = 1, z_{[0:k-1]}$, the distributions of $z_{[0:k-1]}$ are identical. We let P_z denote this distribution. Then, using the separability property of the distance measure over observable $z_{[0:k-1]}$ in the shuffle-then-randomize model, we have:

$$\begin{aligned} & D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X') \mid i_{O_v}(k) = 1) \\ & = \mathbb{E}_{z \sim P_z} D(\mathcal{P}_{s-r}(X) \parallel \mathcal{P}_{s-r}(X') \mid i_{O_v}(k) = 1, z_{[0:k-1]} = z) \\ & \geq \mathbb{E}_{z \sim P_z} D(\mathcal{R}_k(z_{[0:k-1]}, a) \parallel \mathcal{R}_k(z_{[0:k-1]}, b)) \\ & \geq \min_z D(\mathcal{R}_k(z_{[0:k-1]}, a) \parallel \mathcal{R}_k(z_{[0:k-1]}, b)). \end{aligned}$$

APPENDIX D

THEOREM 5 ON SUBSAMPLE-RANDOMIZE-SHUFFLE MODEL UNDER IN-OUT ATTACKS

Theorem 5 (*Subsample-randomize-shuffle model under in-out attack*): Given a protocol \mathcal{P}_{s-r} and sampled user subsets $\{U_k\}_{k \in [K]}$ in the subsample-randomize-shuffle model, and two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$. Let i_{O_v} denote

the in-out information about user i that $i_{o_v}(k) = 1$, and let $S = \{X(i)\}_{i \in U_k}$, $S' = \{X'(i)\}_{i \in U_k} \in \mathbb{X}^{|U_k|}$ denote neighboring datasets w.r.t. U_k , then for any distance measure D that satisfies the data processing inequality:

$$\begin{aligned} & D(\mathcal{P}_{s-r-s}(X) \parallel \mathcal{P}_{s-r-s}(X') \mid U_k, i_{o_v}(k) = 1, z_{(0:k-1)}) \\ & \geq D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0 = z_{(0:k-1)}), \end{aligned}$$

where $\mathcal{R}_{[|U_k|]} = \mathcal{R}_{(k)}$ are the local randomizers of a shuffle-then-randomize protocol \mathcal{P}_{s-r} and z_0 is the global information in \mathcal{P}_{s-r} .

Proof: To prove that:

$$\begin{aligned} & D(\mathcal{P}_{s-r-s}(X) \parallel \mathcal{P}_{s-r-s}(X') \mid U_k, i_{o_v}(k) = 1, z_{(0:k-1)}) \\ & \geq D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0 = z_{(0:k-1)}), \end{aligned}$$

we consider the output of \mathcal{P}_{s-r-s} with observed/fixed variables $U_k, i_{o_v}(k) = 1$ and $z_{(0:k-1)}$. We then define the following post-processing function:

- (1): Remove the $z_{(k')}$ for $k' \in [k+1 : K]$ from the output list;
- (2): Return $z_0 = z_{(0:k-1)}$ and $z_{(k)}$.

Since $z_{(0:k-1)}$ is fixed, the output distribution of above (2) is equivalent to an algorithm \mathcal{P}_{s-r} in the shuffle-randomize model where global information is $z_{(0:k-1)}$. Therefore, according to the data processing inequality, we have:

$$\begin{aligned} & D(\mathcal{P}_{s-r-s}(X) \parallel \mathcal{P}_{s-r-s}(X') \mid U_k, i_{o_v}(k) = 1, z_{(0:k-1)}) \\ & \geq D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0 = z_{(0:k-1)}). \end{aligned}$$

APPENDIX E

PROOF OF THEOREM 2 ON DIVIDE-RANDOMIZE-SHUFFLE MODEL UNDER MESSAGE-LENGTH ATTACKS

We first prove the privacy amplification lower bound, which indicates the destructive power of message-length information attacks. To prove that $D(\mathcal{P}_{d-r-s}(X) \mid U_{k,l}, len_v = l \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l) \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S'))$, we consider the output of \mathcal{P}_{d-r-s} with observed/fixed variables $U_{k,l}, len_v = l$, and define the following post-processing function:

- (1): Remove the $z_{(k')}$ for $k' \in [k+1 : K]$ from the output;
- (2): Remove $\mathcal{R}_{(k)}(x_{i'})$ from $z_{(k)}$ for all $i' \in U_k \setminus U_{k,l}$ to get $z'_{(k)}$;
- (3): Returns $z_0 = z_{(0:k-1)}$ and $z'_{(k)}$.

When $z_{(0:k-1)}$ is fixed, the output distribution of (2) is equal to an algorithm \mathcal{P}_{s-r} in the shuffle-randomize model where global information is $z_{(0:k-1)}$. Thus, we have:

$$\begin{aligned} & D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)}) \\ & \geq D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0 = z_{(0:k-1)}). \end{aligned}$$

Since X and X' differ only at x_v and $v \in U_k$, in two independent runs: $\mathcal{P}_{d-r-s}(X) \mid U_{k,l}, len_v = l, z_{(0:k-1)}$ and $\mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)}$, the distributions of $z_{(0:k-1)}$ are identical. Let P_{z_0} denote this distribution. Then, using the separability property of distance measure over observable $z_{(0:k-1)}$, we have:

$$\begin{aligned} & D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l) \\ & = \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)} = z_0) \\ & \geq \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0). \end{aligned}$$

We then prove the privacy amplification upper bound, indicating the remaining amplification effects. Considering fixed $z_0 = z_{(0:k-1)}$ and fixed $U_{k,l}$, to prove

$$\begin{aligned} & D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)}) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0), \end{aligned}$$

for the output of \mathcal{P}_{s-r} with local randomizers $\mathcal{R}_{k'} \equiv \mathcal{R}_{(k)}$ and $k' \in [|U_{k,l}|]$, we define the following post-processing function:

- (1): The output of \mathcal{P}_{s-r} given input $z_0, \{x_i\}_{i \in U_k}$ is

$$\{\mathcal{R}_{(k)}(z_{(0:k-1)}, x_{\pi^{-1}(1)}), \dots, \mathcal{R}_{(k)}(z_{(0:k-1)}, x_{\pi^{-1}(|U_k|)})\},$$

where $\pi : U_{k,l} \mapsto [|U_{k,l}|]$ is a uniform-random permutation sampled by \mathcal{P}_{s-r} . Now initialize $z_{(k)}$ as the output of \mathcal{P}_{s-r} , for every $i' \in U_k \setminus U_{k,l}$, compute $\mathcal{R}_{(k)}(z_{(0:k-1)}, x_{i'})$, append it to $z_{(k)}$. Then, uniform-randomly permutes the $z_{(k)}$.

- (2): Compute $z_{(k')} = \mathcal{S}(\{\mathcal{R}_{(k')}(z_{(0:k'-1)}, x_{i'})\}_{i' \in U_{k'}})$ for $k' \in [k+1 : K]$ sequentially.
- (3): Return $z_{(0)}, z_{(1)}, \dots, z_{(K)}$.

The output distributions of (3) with $x_v = a$ or $x_v = b$ are equal to the output distributions of $\mathcal{P}_{d-r-s}(X) \mid U_{k,l}, len_v = l, z_{(0:k-1)}$ and $\mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)}$, respectively. According to the data processing inequality and the separability property of distance measure, we have:

$$\begin{aligned} & D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l) \\ & = \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{d-r-s}(X) \parallel \mathcal{P}_{d-r-s}(X') \mid U_{k,l}, len_v = l, z_{(0:k-1)} = z_0) \\ & \leq \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0) \\ & \leq \max_{z_0} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0). \end{aligned}$$

APPENDIX F

JOINT IN-OUT AND MESSAGE LENGTH ATTACKS ON MULTINOMIAL-RANDOMIZE-SHUFFLE MODEL

Theorem 6 (Privacy amplification of MRS model under in-out and message-length attacks): Given a protocol \mathcal{P}_{m-r-s} and sampled non-overlapping and complete user divisions $\{U_k\}_{k \in [K]}$ in the multinomial-randomize-shuffle model, and two neighboring datasets $X = \{x_0, \dots, x_v = a, \dots, x_n\}$, $X' = \{x_0, \dots, x_v = b, \dots, x_n\} \in \mathbb{X}^n$ that differ at the v -th user data. Assuming that $v \in U_k$ and $len(\mathcal{R}_{(k)}(z_{(0:k-1)}, a)) \stackrel{d}{=} len(\mathcal{R}_{(k)}(z_{(0:k-1)}, b))$, let l denote the observed message-length information about user v . Define $U_{k,l}$ as the set of users having the message length (i.e., $U_{k,l} = \{i \mid \text{for } i \in U_k \text{ and } len(\mathcal{R}_{(k)}(x_i)) = l\}$), and let $S = \{X(i)\}_{i \in U_{k,l}}$, $S' = \{X'(i)\}_{i \in U_{k,l}} \in \mathbb{X}^{|U_{k,l}|}$ denote neighboring datasets w.r.t. $U_{k,l}$, then for any distance measure D that satisfies the data processing inequality and the separability property:

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \parallel \mathcal{P}_{m-r-s}(X') \mid U_k, U_{k,l}, len_v = l) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \parallel \mathcal{P}_{s-r}(S') \mid z_0), \end{aligned}$$

where $\mathcal{R}_{[U_{k,l}]} = \mathcal{R}_{(k)}$ are the local randomizers of a shuffle-then-randomize protocol \mathcal{P}_{s-r} and z_0 is the global information in \mathcal{P}_{s-r} .

Proof: First, we prove the privacy amplification lower bound, which highlights the destructive power of message-length information attacks. To prove that:

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | U_k, U_{k,l}, len_v = l) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \|\mathcal{P}_{s-r}(S')), \end{aligned}$$

we consider the output of \mathcal{P}_{m-r-s} with observed/fixed variables $U_k, U_{k,l}, len_v = l$, and define the following post-processing function:

- (1): Remove the $z_{(k')}$ for $k' \in [k+1 : K]$ from the output;
- (2): Remove $\mathcal{R}_{(k)}(x_i)$ from $z_{(k)}$ for all $i \in U_k \setminus U_{k,l}$ to get $z'_{(k)}$;
- (3): Return $z_0 = z_{(0:k-1)}$ and $z'_{(k)}$.

When $z_{(0:k-1)}$ is fixed, the output distribution of (2) equals to an algorithm \mathcal{P}_{s-r} in the shuffle-randomize model where global information is $z_{(0:k-1)}$. Consequently, we obtain

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | U_k, U_{k,l}, len_v = l, z_{(0:k-1)}) \\ & \geq D(\mathcal{P}_{s-r}(S) \|\mathcal{P}_{s-r}(S') | z_0 = z_{(0:k-1)}). \end{aligned}$$

Since X and X' differ only at x_v and $v \in U_k$, in two independent runs: $\mathcal{P}_{d-r-s}(X) | U_k, U_{k,l}, len_v = l, z_{(0:k-1)}$ and $\mathcal{P}_{d-r-s}(X') | U_k, U_{k,l}, len_v = l, z_{(0:k-1)}$, the distributions of $z_{(0:k-1)}$ are identical. We denote this distribution as P_{z_0} . Then, using the separability property of distance measure, we obtain:

$$\begin{aligned} & D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | U_k, U_{k,l}, len_v = l) \\ & = \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | U_k, U_{k,l}, len_v = l, z_{(0:k-1)} = z_0) \\ & \geq \mathbb{E}_{z_0 \sim P_{z_0}} D(\mathcal{P}_{s-r}(S) \|\mathcal{P}_{s-r}(S') | z_0) \\ & \geq \min_{z_0} D(\mathcal{P}_{s-r}(S) \|\mathcal{P}_{s-r}(S') | z_0). \end{aligned}$$

We note that Theorem 6 considers simplified conditional cases with known $U_k, U_{k,l}, len(\mathcal{R}_{(k)}(x_v)) = l$. For unconditional cases aiming at analyze the divergence:

$$D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | len(\mathcal{R}_{(k)}(x_v = b))),$$

if $len(\mathcal{R}_{(k)}(x_v = a))$ follows a different distribution as $len(\mathcal{R}_{(k)}(x_v = b))$, then there will be additional local privacy loss (described earlier); if $len(\mathcal{R}_{(k)}(x_v = a))$ follows the same distribution as $len(\mathcal{R}_{(k)}(x_v = b)) = l$, then for distance measures satisfying the linearity property, since the $\mathbb{P}[U_k, U_{k,l}, len_v = l, z_{(0:k-1)}]$ are identical in two independent runs: $\mathcal{P}_{m-r-s}(X) | U_k, U_{k,l}, len_v = l, z_{(0:k-1)}$ and $\mathcal{P}_{m-r-s}(X') | U_k, U_{k,l}, len_v = l, z_{(0:k-1)}$, the overall divergence can be upper bounded by an expectation of the formulas in Theorem 6, according to the separability property (for observable variable $len_v = l, z_{(0:k-1)}$) and the conditioning increasing property (for unobserved prior distribution of $U_{k,l}$) of the distance measure.

APPENDIX G

PROOF OF THEOREM 3 ON RECTIFIED BIN-RANDOMIZE-SHUFFLE MODEL

Since $len(\mathcal{R}(x)) \equiv len_{max}$ in Algorithm 6, it suffices to prove:

$$\begin{aligned} & D(\mathcal{P}_{b-r-s}(X) \|\mathcal{P}_{b-r-s}(X') | U_{(m)}) \\ & \leq \max_{z_0} D(\mathcal{P}_{m-r-s}(S) \|\mathcal{P}_{m-r-s}(S') | z_0). \end{aligned}$$

Considering fixed $z_0 = z_{(0:k-1)}$, we define the following post-processing function for the output of \mathcal{P}_{m-r-s} with local randomizers $\{\mathcal{R}_{(k')}\}_{k' \in Q_m}$ and query selection distribution $P_{(l)}$:

- (1): Let $z_{(b_{m-1})}, \dots, z_{(b_m-1)}$ denote the output from \mathcal{P}_{m-r-s} with $|Q_m|$ rounds.
- (2): Compute $z_{(k')} = \mathcal{S}_{(k')}(\{\mathcal{R}_{(k')}(z_{(0:k'-1)}), x_i\}_{i \in U_{k'}})$ for $k' \in [b_m : K]$ sequentially, where $U_{k'}$ is chosen by users with the binned multinomial participation paradigm.
- (3): Return $z_{(0)}, z_{(1)}, \dots, z_{(K)}$.

The output distributions of the post-processing with $x_v = a$ or $x_v = b$ are equal to the output distributions of $\mathcal{P}_{b-r-s}(X) | U_{(m)}, z_{(0:k-1)}$ and $\mathcal{P}_{b-r-s}(X') | U_{(m)}, z_{(0:k-1)}$, respectively. According to the data processing inequality and the separability property of distance measure, we have:

$$\begin{aligned} & D(\mathcal{P}_{b-r-s}(X) | U_{(m)} \|\mathcal{P}_{b-r-s}(X') | U_{(m)}) \\ & = \mathbb{E}_{z_{(0:k-1)}} D(\mathcal{P}_{b-r-s}(X) \|\mathcal{P}_{b-r-s}(X') | U_{(m)}, z_{(0:k-1)}) \\ & \leq \mathbb{E}_{z_{(0:k-1)}} D(\mathcal{P}_{m-r-s}(S) \|\mathcal{P}_{m-r-s}(S') | z_0 = z_{(0:k-1)}) \\ & \leq \max_{z_0} D(\mathcal{P}_{m-r-s}(S) \|\mathcal{P}_{m-r-s}(S') | z_0). \end{aligned}$$

APPENDIX H

COMPLEMENTARY RESULTS ABOUT UNCONDITIONAL DIVERGENCES UNDER SIDE-CHANNEL ATTACKS

We note that Theorem 1 considers conditional cases with $io_v(k) = 1$. For the unconditional case $D(\mathcal{P}_{s-r}(X) \|\mathcal{P}_{s-r}(X') | io_v)$, the fact that random variables io_v are distributed identically, regardless of whether X or X' is provided as input, allows for the derivation of a lower bound using the conditioning increasing property of D . Specifically, we have:

$$\begin{aligned} & D(\mathcal{P}_{s-r}(X) \|\mathcal{P}_{s-r}(X') | io_v) \\ & \geq \mathbb{E}_{k \sim \text{uniform}[n]} D(\mathcal{P}_{s-r}(X) \|\mathcal{P}_{s-r}(X') | io_v(k) = 1) \\ & \geq \mathbb{E}_{k \sim \text{uniform}[n]} \min_z D(\mathcal{R}_k(z_{[0:k-1]}, a) \|\mathcal{R}_k(z_{[0:k-1]}, b)). \end{aligned}$$

Theorem 5 considers conditional cases where $U_k, io_v(k) = 1$, and $z_{(0:k-1)}$ is fixed. For unconditional U_k, io_v , the conditioning increasing property of distance measure D can be applied similarly. It is important to note that $z_{(0:k-1)}$ is observable in the subsample-randomize-shuffle model, hence $z_{(0:k-1)}$ consistently appears as a condition.

Theorem 2 considers simplified conditional cases with known $U_{k,l}, len_v = l$. Regarding the unconditional case:

$$D(\mathcal{P}_{m-r-s}(X) \|\mathcal{P}_{m-r-s}(X') | len(\mathcal{R}_{(k)}(b))),$$

if $len(\mathcal{R}_{(k)}(a))$ follows a different distribution as $len(\mathcal{R}_{(k)}(b))$, there is an additional local privacy loss (as described in the previous paragraph). If $len(\mathcal{R}_{(k)}(x_v = a))$ follows the same distribution as $len(\mathcal{R}_{(k)}(x_v = b))$, then for distance measures satisfying conditional composition, since the probabilities $\mathbb{P}[U_{k,l}, len(\mathcal{R}_{(k)}(x_v) = l), z_{(0:k-1)}]$

are identical in the following two independent runs: $\mathcal{P}_{d-r-s}(X)|_{U_{k,l}, len_v} = l, z_{(0:k-1)}$ and $\mathcal{P}_{d-r-s}(X')|_{U_{k,l}, len_v} = l, z_{(0:k-1)}$, the overall divergence can be upper bounded by an expectation of the formulas presented in the theorem, according to the conditioning increasing property.

APPENDIX I

UNCERTAINTY IN SIDE-CHANNEL INFORMATION

In practical side-channel attacks, the adversary may only observe uncertain/noisy in-out, message-length and message-cardinality information. If we let $Blur$ denote some possibly random function on $\text{SideInfo}(i)_{i \in [n]}$, we can then define a variant of DP in the shuffle model with uncertain side-channel information as:

Definition 14 (DP in the curious shuffle model with uncertain side-channel information): Let $Blur$ be some possibly random function that takes $\text{SideInfo}(i)_{i \in [n]}$ as input, a protocol $\mathcal{P} = (\{\mathcal{R}_i\}_{i \in [n]}, \mathcal{A})$ satisfies (ϵ, δ) -differential privacy in the curious shuffle model w.r.t. $Blur$ iff for all neighboring datasets X and $X' \in \mathbb{X}^n$, the $(\mathcal{S} \circ \mathcal{R}_{[n]}(X), Blur(\text{SideInfo}(i)_{i \in [n]}))$ and $(\mathcal{S} \circ \mathcal{R}_{[n]}(X'), Blur(\text{SideInfo}(i)_{i \in [n]}))$ are (ϵ, δ) -indistinguishable.

Similar to certain cases, when the uncertain message-length and message-cardinality information is data-dependent (depends on the differed elements in neighboring datasets X and X'), then there are local privacy loss (see Definition 13 and the following Definition 15). Depending on the concrete $Blur$ function that models the adversaries' uncertainty, the local privacy loss due to uncertain information can be less severe than the certain cases. Even when there is no local privacy loss, there might be shuffle privacy degradation due to uncertain side-channel information. Their consequences can be analyzed in the same way as the certain cases, using the data processing, separability and conditioning increasing properties of DP (introduced in Section V) by taking into consideration the extra randomness in $Blur$ function.

Definition 15 (Local privacy loss due to uncertain message cardinality): Considering a local randomizer $\mathcal{R} : \mathbb{X} \mapsto \mathbb{Y}^*$, let $num(\mathcal{R}(x))$ denote the number of message of $\mathcal{R}(x)$ for a given input $x \in \mathbb{X}$. Let $Blur$ be some possibly random function that takes $num(\mathcal{R}(x))$ as input, the local privacy loss of uncertain message cardinality information is (ϵ, δ) if following condition holds for some $x, x' \in \mathbb{X}$:

$$D_{e^\epsilon}(Blur(num(\mathcal{R}(x))) \parallel Blur(num(\mathcal{R}(x')))) \geq \delta.$$