

# Privacy-Preserving Data Deduplication for Enhancing Federated Learning of Language Models

Aydin Abadi\*  
Newcastle University  
aydin.abadi@ncl.ac.uk

Vishnu Asutosh Dasu\*  
Pennsylvania State University  
vdasu@psu.edu

Sumanta Sarkar\*  
University of Warwick  
sumanta.sarkar@warwick.ac.uk

**Abstract**—Deduplication is a vital preprocessing step that enhances machine learning model performance and saves training time and energy. However, enhancing federated learning through deduplication poses challenges, especially regarding scalability and potential privacy violations if deduplication involves sharing all clients’ data. In this paper, we address the problem of deduplication in a federated setup by introducing a pioneering protocol, Efficient Privacy-Preserving Multi-Party Deduplication (EP-MPD). It efficiently removes duplicates from multiple clients’ datasets without compromising data privacy. EP-MPD is constructed in a modular fashion, utilizing two novel variants of the Private Set Intersection protocol. Our extensive experiments demonstrate the significant benefits of deduplication in federated learning of large language models. For instance, we observe up to 19.62% improvement in perplexity and up to 27.95% reduction in running time while varying the duplication level between 10% and 30%. EP-MPD effectively balances privacy and performance in federated learning, making it a valuable solution for large-scale applications.

## I. INTRODUCTION

Machine learning (ML) is a data-driven process where the *quality* of the training data significantly influences the *accuracy* of an ML model. Data generated from real-world applications often lacks organization, leading to issues such as missing values, typos, format mismatches, outliers, or duplicated entries within the raw dataset. To ensure meaningful learning, the collected data must undergo a thorough data cleaning process [1]. Duplicated sequences are prevalent in text datasets. They can adversely affect the training process of Language Models (LMs). Lee et al. [2] investigated the Colossal Clean Crawled Corpus (C4) [3] dataset. They discovered a 61-word sequence within C4 that was repeated 61,036 times verbatim in the training dataset and 61 times in the validation set. As a result, the trained model generalized poorly on the dataset, and over 1% of the unprompted model outputs were memorized and copied verbatim from the training dataset.

\*Equal contribution. Listing order is alphabetical.

Upon deduplicating the dataset, they reduced memorization by up to 10× and, in certain cases, even improved perplexity by up to 10%.

Additionally, memorization negatively affects the privacy and fairness of LMs [4]. Memorization makes language models vulnerable to membership inference attacks [5], [6], data extraction attacks [7], [8], and can lead to copyright violations as LMs can regurgitate training data sequences from copyrighted sources [9]. Carlini et al. [4] conclude that bigger LMs memorize more and more duplicates increase memorization. Kandpal et al. [10] show that the success of most privacy attacks on LMs is largely due to the duplication in the training datasets. Furthermore, in the FL setting, malicious clients can exploit the memorization of LMs to extract sensitive information from honest clients’ datasets [11]. Therefore, we must adopt data-cleaning practices to improve model utility and reduce the risks of privacy attacks.

While deduplication can improve model performance and reduce memorization, it also enhances training efficiency in various aspects [2]. Removing duplicates reduces GPU time and minimizes training costs in terms of time, money, and carbon footprint. This streamlined process optimizes resource utilization and contributes to more sustainable ML practices [12], [13], [14]. Sustainable development is a global priority, and aligning ML practices with this theme is crucial.

Federated learning (FL) is collaborative learning that allows training across multiple decentralized devices without exchanging data. In FL, devices compute local models based on their data and then share the local model updates with a central server. This server aggregates the updates to derive a global model that encapsulates the features of all the local data held by the individual devices [15]. Several real-world applications involve training with FL such as healthcare [16], smart city [17], and edge computing [18]. Some of these applications pose a risk of duplicates in the local training data across multiple devices. For example, Google Keyboard suggestions from a user’s text query rely on FL [19]. Local models are trained in situ on Android phones with user data. In this setting, many text queries typed by the users across multiple phones are the same.

For an FL process to be efficient and effective, participating devices must perform deduplication. When a data owner solely

intends to remove duplicates from their dataset, the privacy risk is minimal since the data is entirely under the owner’s control. In the context of FL, the scenario shifts significantly when deduplication is introduced. Multiple devices participating in FL may possess overlapping data, even after deduplicating their datasets. If they aim to deduplicate the combined dataset, one approach could involve sharing their raw data with each other and checking for intersections. However, this approach compromises privacy.

Hence, there is a necessity for privacy-preserving deduplication in FL, where participating devices would collaboratively deduplicate the combined datasets in a privacy-preserving fashion. This paper demonstrates how deduplication can be executed in FL without compromising the privacy of the data belonging to the individual devices.

### A. Overview of Privacy-Preserving Deduplication in FL

Consider nodes  $D_1, \dots, D_n$  that are involved in FL, where each  $D_i$  has a dataset  $S_i$ . Effectively, the training in FL is conducted on the union of these datasets, i.e.,  $\bigcup_{i=1}^n S_i$ . If intersections exist among the datasets, FL will incorporate duplicates, leading to the drawbacks discussed earlier. Consider nodes  $D_1$  and  $D_2$ , each with datasets  $S_1$  and  $S_2$ , respectively implying that FL training occurs on  $S_1 \cup S_2$ . If the intersection  $S_1 \cap S_2$  is nonempty, training individually on  $S_1$  and  $S_2$  would mean training twice on the duplicate  $S_1 \cap S_2$ . Thus, one of  $D_1$  and  $D_2$  should remove the duplicate  $S_1 \cap S_2$ . This should be done without harming each other’s data privacy. Our solution applies private set intersection (PSI) which securely finds the intersection of  $S_1$  and  $S_2$  without revealing any other elements [20], [21], [22]. Once  $D_1$  and  $D_2$  learn  $S_1 \cap S_2$  through PSI,  $D_1$  will train on  $S'_1 = S_1 \setminus S_1 \cap S_2$  and  $D_2$  will train on  $S_2$ . According to set theory, it holds that  $S'_1 \cup S_2 = S_1 \cup S_2$ . Hence, the resulting FL model remains as intended, while the training process is devoid of duplicates, avoiding the associated drawbacks. The scenario with two nodes seems straightforward. However, it becomes complicated when more than two nodes participate. We now consider  $n$  nodes  $D_1, \dots, D_n$ , where each node  $D_i$  has a set  $S_i, \forall i, 1 \leq i \leq n$ , and  $n > 2$ . Following the approach used in the two-node case, one might be tempted to (i) apply multi-party PSI to  $n$  nodes, (ii) find their intersection  $I_n$ , and (iii) let node  $D_1$  train on  $S'_1 = S_1 \setminus I_n$ , and rest of the nodes train on their own dataset  $S_i$ . This method removes the duplicates that exist across *all* the nodes. However, this does not detect and remove the duplicates that exist among a subset of nodes. For instance, if there is a subset of  $k$  nodes ( $k < n$ ) with a large intersection  $I_k$ , then  $I_k$  will remain in the full training dataset  $S'_1 \cup (\bigcup_{i=2}^n S_i)$  as duplicates. A generic multi-party PSI does not help in this case. Therefore, we need to consider each pair of nodes and remove the duplicates accordingly.

### B. Our Contributions

In this paper, our end goal is to develop a scheme that allows multiple clients to benefit from deduplication while training on

their data in a federated learning setup. The first requirement is an efficient deduplication of sets belonging to multiple clients. To address this, we introduce the notion of *Privacy-Preserving Multi-Party Deduplication* (P-MPD) in Section IV. This essentially describes a functionality that takes input datasets  $S_1, \dots, S_m$  (potentially containing duplicates) from  $m$  clients and outputs the datasets  $S'_1, \dots, S'_m$  such that  $\bigcup_{i=1}^m S'_i = \bigcup_{i=1}^m S_i$ , where  $S'_i \cap S'_j = \emptyset, i \neq j$ . We realize that an efficient construction of P-MPD requires a substantially improved PSI protocol that supports scalability. This leads us to introduce a new notion for PSI which we call *Group PSI* (G-PSI) in Section III. The functionality of G-PSI takes sets from a group of clients and returns each client the intersection of their sets with all the other clients’ sets. We provide constructions of *Efficient Group PSI* (EG-PSI) that efficiently realizes G-PSI, namely EG-PSI<sup>(1)</sup> and EG-PSI<sup>(II)</sup> in Figures 1 and 2, respectively. EG-PSI<sup>(1)</sup> is based on symmetric key primitives such as pseudorandom permutation, while EG-PSI<sup>(II)</sup> is developed using an oblivious pseudorandom function, which is a well-known public key primitive. With the building block EG-PSI, we build *Efficient Privacy-Preserving Multi-Party Deduplication* (EP-MPD) that realizes P-MPD as shown in Figure 4. We prove the security of EP-MPD, EG-PSI<sup>(1)</sup>, and EG-PSI<sup>(II)</sup> within the simulation-based paradigm. Our construction of EP-MPD allows for efficient duplicate removal without compromising clients’ data privacy, resulting in an improved model after running federated learning on the deduplicated datasets, as outlined in Figure 5.

We perform an extensive experimental evaluation to benchmark EP-MPD and the effect of the resulting deduplication on FL. The overall running time for EP-MPD<sup>(1)</sup>, which employs EG-PSI<sup>(1)</sup> (symmetric key primitives based), is much less than EP-MPD<sup>(II)</sup>, which utilizes EG-PSI<sup>(II)</sup> (public key primitives based). Overall, clients enjoy relatively less computation time in EP-MPD<sup>(II)</sup>. Our protocols can scale to large datasets and client counts. For example, when 50 clients have  $2^{19}$  data points in their datasets comprising of 30% duplicates, then EP-MPD<sup>(1)</sup> takes 1160 seconds and EP-MPD<sup>(II)</sup> takes 7653 seconds; client running time is 641 and 111 seconds in EP-MPD<sup>(1)</sup> and EP-MPD<sup>(II)</sup> respectively. We experiment with fine-tuning 2 LMs with 7 datasets and 10 clients in FL. We achieve an improvement of up to 19.62% in perplexity and up to 27.95% improvement in GPU training time while varying the duplication level between 10% and 30%.

## II. PRELIMINARIES

### A. Notations and Assumptions

We define a wrapper function  $\text{Update}(S, \hat{S}) \rightarrow S$  which takes two sets,  $S$  and  $\hat{S}$ . It updates  $S$  by removing from it the elements in set  $\hat{S}$  and returns the updated set  $S$ . In this paper, by the sum of sets (e.g.,  $\sum_{i=1}^n S_i$ ) we mean the concatenation of the sets which may result in a multi-set. We denote an empty set by  $\emptyset$ . We denote a size of vector  $\vec{v}$  with  $|\vec{v}|$ . We assume that the server and all the users have access to secure channels among them. By the notation,  $\mathcal{X} \stackrel{c}{=} \mathcal{Y}$ , we mean that the two distributions  $\mathcal{X}$  and  $\mathcal{Y}$  are computationally indistinguishable.

## B. Federated Learning (FL)

The concept of FL was proposed by McMahan et al. [15] as a framework for training an ML model where the training data are distributed across multiple devices. In FL, a server orchestrates the training of a global model by aggregating models locally computed by the clients on their local devices. Suppose there are  $n$  clients, and each client  $D_i$  has a dataset  $S_i$ . The server has the initial model  $\theta$ . It sends  $\theta$  to the clients, and each client performs gradient descent computation on their local dataset  $S_i$  as  $J_i(S_i, \theta) = \frac{1}{d_i} \sum_{(\mathbf{x}, y) \in S_i} C(\theta, (\mathbf{x}, y))$ , where  $C$  is the cost function and  $d_i = |S_i|$ . The client  $D_i$  computes the local models as  $\theta_i \leftarrow \theta - \eta \nabla J(S_i, \theta)$ , where  $\eta$  is the learning rate. After receiving the local models from the clients, the server performs an aggregated averaging on the local models to derive the global model as:

$$\Theta = \frac{1}{d} (d_1 \theta_1 + \dots + d_n \theta_n) \quad (1)$$

where  $d = \sum_{i=1}^n d_i$  is the total size of the dataset  $S = \sum_{i=1}^n S_i$ , and  $\theta_i$  is the locally trained model on the dataset  $S_i$ . This computation is repeated until the model converges.

While this framework achieves a basic level of privacy in which each user’s data is not directly sent to the server, it is susceptible to advanced privacy attacks such as membership inference attacks [5], [6] and data extraction attacks [8], [11], [23], [24] as the client models  $\theta_i$  are aggregated in a non-private fashion on the server side. Some privacy-preserving FL protocols have attempted to mitigate these attacks by securely aggregating the client models. These protocols rely on cryptographic techniques like homomorphic encryption [25], [26], functional encryption [27], or secure aggregation [28], [29]. Additionally, differentially private training techniques [30] can be used to ensure differential privacy guarantees on the global model  $\Theta$ . We emphasize that in this paper, our proposed schemes are agnostic to the type of FL mechanism used.

## C. Causal Language Modeling (CLM)

Causal Language Modeling (CLM) is a natural language processing task where the goal of the language model is to predict the next word or token given a sequence of tokens. The language model autoregressively generates the next token until a pre-determined sequence length is reached or a special STOP token is generated. Given a sequence of  $n$  tokens  $\mathbf{Y} = \{y_1, y_2, \dots, y_{n-1}, y_n\}$ , the language model  $\Theta$  is trained to learn the following probability distribution:

$$P(\mathbf{Y}) = \prod_{i=1}^n P(y_i | y_1, \dots, y_{i-1}) \quad (2)$$

The CLM training objective is to minimize the negative log-likelihood loss given by:

$$\mathcal{L}(\Theta, \mathbf{Y}) = - \sum_{i=1}^n \log(\Theta(y_i | y_1, \dots, y_{i-1})) \quad (3)$$

After training is complete, the text is autoregressively sampled from the language model i.e.,  $\hat{y}_{i < n} \sim \Theta(y_i | y_1, \dots, y_{i-1})$ .

The *perplexity* metric is commonly used to evaluate the performance of the language model to determine how well it has learned the probability distribution in Equation 2. The perplexity  $PP$  of a sequence  $\mathbf{y}$  is defined as:

$$PP(\mathbf{Y}) = \exp \left( -\frac{1}{n} \sum_{i=1}^n \log(\Theta(y_i | y_1, \dots, y_{i-1})) \right) \quad (4)$$

A lower perplexity score implies that model  $\Theta$  has been trained well to estimate the real-world probability distribution. Informally, a sequence with a low perplexity score implies that the model is less “surprised” by a sequence of tokens.

## D. Security Model

In this paper, we use the simulation-based paradigm of secure multi-party computation [31] to define and prove the proposed protocol. Since we focus on the static passive (semi-honest) adversarial model, we will restate the security definition within this context, after outlining the threat model.

1) *Threat Model Outline*: In this paper, three types of parties are involved: clients, a server, and a third party. We allow these parties to be corrupted by semi-honest adversaries. In this adversarial model, parties follow the protocols’ instructions. Therefore, adversaries do not modify the model architecture to better suit their attack or send malicious (global) messages or model parameters. They may try to learn more information (from the messages they exchange) than they are supposed to learn [31]. In the protocols, we instantiate the third party using a trusted execution environment ( $\mathcal{T}\mathcal{E}\mathcal{E}$ ).  $\mathcal{T}\mathcal{E}\mathcal{E}$ s are commonly used to instantiate a trusted third party due to their presumed security guarantees. The core assurance about  $\mathcal{T}\mathcal{E}\mathcal{E}$ s is that they are fully trusted. However, recent studies have shown that  $\mathcal{T}\mathcal{E}\mathcal{E}$ s like Intel SGX are vulnerable [32]. Therefore, in our solution, we do not assume  $\mathcal{T}\mathcal{E}\mathcal{E}$ s to be fully trusted and minimize our trust assumption to being semi-honest<sup>1</sup>. Specifically, we assume that  $\mathcal{T}\mathcal{E}\mathcal{E}$  does not collude with other parties. Furthermore, we allow any  $n - 1$  clients to collude with each other (where  $n$  is the total number of clients).

2) *Two-party Computation*: A two-party protocol  $\Gamma$  is captured by specifying a random process that maps a pair of inputs to a pair of outputs (one output for each party). Such process is referred to as a functionality denoted by  $f : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^* \times \{0, 1\}^*$ , where  $f := (f_1, f_2)$ . For every input pair  $(x, y)$ , the output pair is a random variable  $(f_1(x, y), f_2(x, y))$ , such that the party with input  $x$  wishes to obtain  $f_1(x, y)$  while the party with input  $y$  wishes to receive  $f_2(x, y)$ . The above functionality can be easily extended to more than two parties.

3) *Security in the Presence of Passive Adversaries*: In the passive adversarial model, the party corrupted by such an adversary correctly follows the protocol specification. Loosely speaking, a protocol is secure if whatever can be computed by a corrupt party in the protocol can be computed using its input

<sup>1</sup>A similar threat model for the semi-honest  $\mathcal{T}\mathcal{E}\mathcal{E}$  can be found in [33].

and output only. In the simulation-based model, it is required that a party's "view" in a protocol's execution can be simulated given only its input and output. This implies that the parties learn nothing from the protocol's execution.

Formally, in two-party case, party  $i$ 's view (during the execution of  $\Gamma$ ) on input pair  $(x, y)$  is denoted by  $\text{View}_i^\Gamma(x, y)$  and equals  $(w, r^i, m_1^i, \dots, m_i^i)$ , where  $w \in \{x, y\}$  is the input of the  $i^{\text{th}}$  party,  $r_i$  is the outcome of this party's internal random coin tosses, and  $m_j^i$  represents the  $j^{\text{th}}$  message this party receives. The output of the  $i^{\text{th}}$  party during the execution of  $\Gamma$  on  $(x, y)$  is denoted by  $\text{Output}_i^\Gamma(x, y)$  and can be generated from its own view of the execution.

**Definition 1.** Let  $f$  be the deterministic functionality defined above. Protocol  $\Gamma$  securely computes  $f$  in the presence of a static passive probabilistic polynomial-time (PPT) adversary  $\mathcal{A}$ , if for every  $\mathcal{A}$  in the real model, there exist PPT algorithms  $(\text{Sim}_1, \text{Sim}_2)$  such that:

$$\{\text{Sim}_1(x, f_1(x, y))\}_{x, y} \stackrel{c}{=} \{\text{View}_1^{\mathcal{A}, \Gamma}(x, y)\}_{x, y}$$

$$\{\text{Sim}_2(y, f_2(x, y))\}_{x, y} \stackrel{c}{=} \{\text{View}_2^{\mathcal{A}, \Gamma}(x, y)\}_{x, y}$$

Definition 1 can be easily extended to  $m > 2$  parties.

#### E. Pseudorandom Function and Permutation

Informally, a pseudorandom function  $\text{PRF}(\cdot)$  is a deterministic function that takes a key of length  $\lambda$  and an input of length  $u$ ; and outputs a value of length  $v$  indistinguishable from an output of a truly random function. More formally, a pseudorandom function can be defined as  $\text{PRF} : \{0, 1\}^\lambda \times \{0, 1\}^u \rightarrow \{0, 1\}^v$ , where  $\lambda$  is the security parameter.

The definition of a pseudorandom permutation,  $\text{PRP} : \{0, 1\}^\lambda \times \{0, 1\}^u \rightarrow \{0, 1\}^u$ , is very similar to that of a pseudorandom function, with a difference; namely, it is required the keyed function  $\text{PRP}(k, \cdot)$  to be indistinguishable from a uniform permutation, instead of a uniform function. In cryptographic schemes that involve PRP, sometimes honest parties may be required to compute the inverse of pseudorandom permutation, i.e.,  $\text{PRP}^{-1}(k, \cdot)$ , as well. In this case, it would require that  $\text{PRP}(k, \cdot)$  be indistinguishable from a uniform permutation even if the distinguisher is additionally given oracle access to the inverse of the permutation.

#### F. Oblivious Pseudorandom Function

An Oblivious Pseudorandom Function (OPRF) is a protocol that involves a client and a server. OPRF enables the client with input  $x \in \{0, 1\}^u$ , and the server with key  $k \in \{0, 1\}^\lambda$  to execute an instance of PRF. Informally, the security of an OPRF asserts that, by the completion of OPRF, the client only learns the output of the PRF evaluated on inputs  $k$  and  $x$ , i.e.,  $\text{PRF}(k, x)$  while the server gains no information, e.g., about the input of the client and the output of PRF.

#### G. Trusted Execution Environments

Trusted Execution Environment ( $\mathcal{T}\mathcal{E}\mathcal{E}$ ), also known as a secure enclave, constitutes a secure processing environment comprising processing, memory, and storage hardware units

[34], [35]. Within this environment, the code and data residing in them remain isolated from other layers in the software stack, including the operating system. An ideal  $\mathcal{T}\mathcal{E}\mathcal{E}$  guarantees the preservation of data integrity and confidentiality. Given the assumption that the physical CPU remains uncompromised, enclaves are shielded from attackers with physical access to the machine, including the memory and system bus. Side-channel attacks on different deployments of  $\mathcal{T}\mathcal{E}\mathcal{E}$ s have been demonstrated in the literature [36]. These attacks pose a threat as they could enable attackers to extract secrets from  $\mathcal{T}\mathcal{E}\mathcal{E}$ s. Nevertheless,  $\mathcal{T}\mathcal{E}\mathcal{E}$ s technologies have been evolving to address and mitigate side-channel attacks.

Our security, trust, and system assumptions regarding  $\mathcal{T}\mathcal{E}\mathcal{E}$ s are conservative. Specifically, as detailed in Sections IV-B1 and III-B, our solution (i) avoids disclosing any plaintext messages or private keys to  $\mathcal{T}\mathcal{E}\mathcal{E}$  and (ii) does not expect  $\mathcal{T}\mathcal{E}\mathcal{E}$  to maintain an extensive storage and memory space or possess strong processing resources. Instead, we establish a formal model and construction under the assumption that  $\mathcal{T}\mathcal{E}\mathcal{E}$  guarantees execution integrity (and authenticity), and ensures minimal confidentiality. Specifically, we formally demonstrate that  $\mathcal{T}\mathcal{E}\mathcal{E}$  at the most only learns the size of the encrypted computation result (i.e., the intersections' cardinality). In our work,  $\mathcal{T}\mathcal{E}\mathcal{E}$  can be seamlessly substituted with any semi-honest server that does not collude with other entities.

### III. GROUP PSI (G-PSI)

#### A. Formal Definition of G-PSI

In this section, we present the concept of Group PSI (G-PSI). In G-PSI, there are two groups of clients,  $\mathcal{G}_0$  and  $\mathcal{G}_1$ , where each group contains  $m$  clients,  $\mathcal{G}_j : \{\mathcal{C}_{j,1}, \dots, \mathcal{C}_{j,m}\}$ ,  $0 \leq j \leq 1$ . Each client  $\mathcal{C}_{j,1}$  has a set  $S_{j,1}$ , such that no pair of clients' sets in the same group share a common element.

Informally, G-PSI allows every client in one group to (efficiently) find the intersection that their set has with the set of every client of the other group, without allowing them to learn anything beyond that about other clients' set elements. To achieve high computational efficiency in G-PSI, we will involve a third-party  $\mathcal{T}\mathcal{P}$  that assists the clients with computing the intersection. The functionality  $f_{\text{G-PSI}}$  that G-PSI computes takes a set  $S_{j,i}$  from every client  $\mathcal{C}_{j,i}$  and no input from  $\mathcal{T}\mathcal{P}$ . It returns (i) to every client  $\mathcal{C}_{j,i}$  a vector  $\vec{v}_{j,i}$  which contains the intersection that  $\mathcal{C}_{j,i}$ 's set has with every other client's set in the other group and (ii) to  $\mathcal{T}\mathcal{P}$  an empty set  $\emptyset$ . Hence, functionality  $f_{\text{G-PSI}}$  can be formally defined as follows.

$$f_{\text{G-PSI}}\left(\underbrace{(S_{0,1}, \dots, S_{0,m})}_{\mathcal{G}_0}, \underbrace{(S_{1,1}, \dots, S_{1,m})}_{\mathcal{G}_1}, \emptyset\right) \rightarrow \left(\underbrace{(\vec{v}_{0,1}, \dots, \vec{v}_{0,m})}_{\mathcal{G}_0}, \underbrace{(\vec{v}_{1,1}, \dots, \vec{v}_{1,m})}_{\mathcal{G}_1}, \emptyset\right), \quad (5)$$

where,  $\vec{v}_{j,i} = [S_{j,i} \cap S_{1-j,1}, \dots, S_{j,i} \cap S_{1-j,m}]$ ,  $0 \leq j \leq 1$ , and  $1 \leq i \leq m$ .

In the case where clients of only one of the groups, e.g.,  $\mathcal{G}_0$ , receives the result, then the above functionality simply returns  $\emptyset$  to the clients in the other group, e.g.,  $\vec{v}_{1,1} = \dots = \vec{v}_{1,m} = \emptyset$ .

Since  $\mathcal{TP}$  performs computation on all clients' encrypted sets, there is a possibility of leakage to  $\mathcal{TP}$ . Depending on the protocol that realizes  $f_{\text{G-PSI}}$  this leakage could contain different types of information; for instance, it could contain (i) the size of the intersection of any two clients' sets, or (ii) the size of the intersection of all clients' sets, or (iii) nothing at all. Often such leakage is defined by a leakage function  $\mathcal{L}$  that takes all parties (encoded) inputs and returns the amount of leakage.

We assert that a protocol securely realizes  $f_{\text{G-PSI}}$  if (1) it reveals nothing beyond a predefined leakage to  $\mathcal{TP}$  and (2) whatever can be computed by a client in the protocol can be obtained from its input and output. This is formalized under the simulation model. We require a client's view during G-PSI's execution to be simulatable given its input and output. We require that  $\mathcal{TP}$ 's view can be simulated given the leakage.

**Definition 2** (Security of G-PSI). Let  $\mathcal{G}_0$  and  $\mathcal{G}_1$  be two groups of clients, where each group contains  $m$  clients,  $\mathcal{G}_j : \{\mathcal{C}_{j,1}, \dots, \mathcal{C}_{j,m}\}$ ,  $0 \leq j \leq 1$ . Let  $\mathcal{L}$  denote a leakage function,  $f_{\text{G-PSI}}$  be the functionality defined above (in Relation 5 on page 4),  $S = \{S_{0,1}, \dots, S_{0,m}, S_{1,1}, \dots, S_{1,m}\}$ , and  $S_{j,i}$  represent a set belonging to client  $\mathcal{C}_{j,i}$ . Then, a protocol  $\Gamma$  securely realizes  $f_{\text{G-PSI}}$  in the presence of a static semi-honest PPT adversary  $\mathcal{A}$ , if for every  $\mathcal{A}$  in the real model, there exists a PPT adversary (simulator) Sim in the ideal model, such that for every  $\mathcal{C}_{j,i} \in \{\mathcal{C}_{0,1}, \dots, \mathcal{C}_{0,m}, \mathcal{C}_{1,1}, \dots, \mathcal{C}_{1,m}\}$  and  $\mathcal{TP}$ , Relations 6 and 7 hold respectively.

$$\{\text{Sim}_{\mathcal{C}_{j,i}}(S_{j,i}, \vec{v}_{j,i})\}_S \stackrel{c}{\equiv} \{\text{View}_{\mathcal{C}_{j,i}}^{\mathcal{A}, \Gamma}(S, \emptyset)\}_S \quad (6)$$

$$\{\text{Sim}_{\mathcal{TP}}^c(\emptyset, \emptyset)\}_S \stackrel{c}{\equiv} \{\text{View}_{\mathcal{TP}}^{\mathcal{A}, \Gamma}(S, \emptyset)\}_S \quad (7)$$

where  $\vec{v}_{j,i} = \left[ |S_{j,i} \cap S_{1-j,1}|, \dots, |S_{j,i} \cap S_{1-j,m}| \right]$ ,  $0 \leq j \leq 1$ ,  $1 \leq i \leq m$ , and  $\mathcal{TP}$ 's input is  $\emptyset$ .

### B. Efficient Construction of G-PSI

In this section, we introduce two efficient protocols that realize G-PSI. The first one, called EG-PSI<sup>(l)</sup>, is highly efficient and based on symmetric key cryptography. The second one, called EG-PSI<sup>(ll)</sup>, is based on OPRF (in turn depends on public key cryptography) and discloses less information to  $\mathcal{TEE}$  than the former does. Thus, these two protocols trade-off between performance and leakage amount.

1) EG-PSI<sup>(l)</sup>: At a high level, EG-PSI<sup>(l)</sup> operates as follows. Initially, each client in group  $\mathcal{G}_0$  agrees with every client in group  $\mathcal{G}_1$  on a secret key. Every client encrypts its set elements using every key it has agreed on with other clients (in a different group). Each client, for every encrypted set element, temporally stores a triple that includes (i) the encrypted element, (ii) the key used to encrypt that element, and (iii) the index of the client with whom the key was shared. These triples will enable the client to efficiently (a) retrieve the correct key and (b) identify the client with whom it has the element in common when presented with an encrypted element.

Once all elements are encrypted using the corresponding keys, each client transmits only its encrypted elements to  $\mathcal{TEE}$ .

Given the sets of encrypted elements from all clients,  $\mathcal{TEE}$  aggregates these sets and identifies the encrypted elements that appear more than once. Subsequently,  $\mathcal{TEE}$  forwards to a client those encrypted elements that (1) appear more than once and (2) are among the messages that the client initially sent to  $\mathcal{TEE}$ . Upon receiving each encrypted element from  $\mathcal{TEE}$ , a client searches its local list of triples to locate the corresponding key and the index  $l$  representing a specific client. Utilizing the key, the client decrypts the element and regards the resultant element as one of the elements within the intersection it shares with the  $l$ -th client. For a comprehensive description of EG-PSI<sup>(l)</sup>, refer to Figure 1.

The only information that  $\mathcal{TEE}$  learns in EG-PSI is the *size* of the intersection of any two clients' sets, called *pair-wise intersection cardinality*. Below, we formally define it.

**Definition 3** (pair-wise intersection cardinality). Let  $\underbrace{(S'_{0,1}, \dots, S'_{0,m})}_{\mathcal{G}_0}$ ,  $\underbrace{(S'_{1,1}, \dots, S'_{1,m})}_{\mathcal{G}_1}$  be two groups  $\mathcal{G}_0$  and  $\mathcal{G}_1$  of (encrypted) sets. Then, vector  $\vec{s}$  represents the pair-wise intersection cardinality:  $\vec{s} = [\vec{s}_{0,1}, \dots, \vec{s}_{0,m}, \vec{s}_{1,1}, \dots, \vec{s}_{1,m}]$ , where  $\vec{s}_{j,i} = \left[ |S_{j,i} \cap S_{1-j,1}|, \dots, |S_{j,i} \cap S_{1-j,m}| \right]$ ,  $0 \leq j \leq 1$ , and  $1 \leq i \leq m$ .

**Definition 4.** Let  $\vec{s}$  be a pair-wise intersection cardinality of two groups  $\mathcal{G}_0$  and  $\mathcal{G}_1$  of encrypted sets:  $\underbrace{(S'_{0,1}, \dots, S'_{0,m})}_{\mathcal{G}_0}$ ,  $\underbrace{(S'_{1,1}, \dots, S'_{1,m})}_{\mathcal{G}_1}$ , with respect to Definition 3. Then, leakage function  $\mathcal{L}$  is defined as follows:  $\mathcal{L}\left(\underbrace{(S'_{0,1}, \dots, S'_{0,m})}_{\mathcal{G}_0}, \underbrace{(S'_{1,1}, \dots, S'_{1,m})}_{\mathcal{G}_1}\right) \rightarrow \vec{s}$ .

**Theorem 1.** Let  $f_{\text{G-PSI}}$  be the functionality defined in Relation 5. Let  $\mathcal{L}$  be the leakage function presented in Definition 4. If PRP is a secure pseudorandom permutation, EG-PSI<sup>(l)</sup> (presented in Figure 1) securely realizes  $f_{\text{G-PSI}}$ , w.r.t. Definition 2.

We refer to Appendix A for the proof of Theorem 1.

2) EG-PSI<sup>(ll)</sup>: In this section, we present EG-PSI<sup>(ll)</sup> which is the second variant of EG-PSI. Note that  $\mathcal{TEE}$  in EG-PSI<sup>(l)</sup> is able to learn the cardinality of the intersection of each pair of clients' sets. We aim to reduce this leakage in EG-PSI<sup>(ll)</sup>. However, EG-PSI<sup>(ll)</sup> is no longer based on symmetric key cryptography as it depends on OPRF. To the best of our knowledge, all the constructions of OPRFs are based on public key cryptography whose security depends on some hard problem assumptions. The same holds for EG-PSI<sup>(ll)</sup>.

A brief overview of the construction of EG-PSI<sup>(ll)</sup> is as follows. We have two groups of clients  $\mathcal{G}_0$  and  $\mathcal{G}_1$  having their own set that they want to find the pairwise intersection. There is a  $\mathcal{TEE}$  that has a key  $k$  for a PRF. During the setup, each client interacts with  $\mathcal{TEE}$  to encrypt their set using  $\text{PRF}(k, \cdot)$  through an OPRF call. Then, each client of group  $\mathcal{G}_1$  will send their encrypted set to every client of  $\mathcal{G}_0$ . Upon receiving an encrypted set from a client of  $\mathcal{G}_1$ , each client of  $\mathcal{G}_0$  determines the intersection with their encrypted set. Once the client finds

- *Parties.* Trusted execution environment  $\mathcal{T}\mathcal{E}\mathcal{E}$ , clients in group  $\mathcal{G}_0 : \{\mathcal{C}_{0,1}, \dots, \mathcal{C}_{0,m}\}$ , and clients in group  $\mathcal{G}_1 : \{\mathcal{C}_{1,1}, \dots, \mathcal{C}_{1,m}\}$ .
- *Inputs.* Sets  $S_{0,1}, \dots, S_{0,m}, S_{1,1}, \dots, S_{1,m}$ , where each  $S_{j,i}$  belongs to client  $\mathcal{C}_{j,i}$ ,  $0 \leq j \leq 1$  and  $1 \leq i \leq m$ .
- *Outputs.*  $\vec{v}_{j,i}$  to  $\mathcal{C}_{j,i}$ , where  $\vec{v}_{j,i} = \left[ [S_{j,i} \cap S_{1-j,1}], \dots, [S_{j,i} \cap S_{1-j,m}] \right]$ .

1) Setup.

- a) each client  $\mathcal{C}_{0,i}$  in  $\mathcal{G}_0$  agrees with every client  $\mathcal{C}_{1,l}$  in  $\mathcal{G}_1$  on a secret key  $k_{i,l}$ , by picking a random key  $k_{i,l}$  and sending it to  $\mathcal{C}_{1,l}$ . Client  $\mathcal{C}_{0,i}$  stores this key as  $k_{i,l}$  while  $\mathcal{C}_{1,l}$  stores this key as  $k_{i,l}$ .
- b) each  $\mathcal{C}_{j,i}$  takes the following steps:
  - i) encrypts its set elements under keys  $k_{i,l}$  ( $\forall l, 1 \leq l \leq m$ ) as follows,  $\forall e \in S_{j,i} : \text{PRP}(k_{i,l}, e) \rightarrow e'_{i,l}$ . Let set  $S'_{j,i}$  contain the encrypted set elements of  $\mathcal{C}_{j,i}$  and let set  $T_{j,i}$  contains all triples of the form  $(e'_{i,l}, k_{i,l}, l)$ .
  - ii) sends  $S'_{j,i}$  to  $\mathcal{T}\mathcal{E}\mathcal{E}$  and locally keeps  $T_{j,i}$ .

2) Finding Encrypted Intersection.  $\mathcal{T}\mathcal{E}\mathcal{E}$  takes the following steps for each  $\mathcal{C}_{j,i}$ .

- a) appends to an empty set,  $R_{j,i}$ , every ciphertext that satisfy the following conditions hold:
  - it appears more than once in the set  $S = \sum_{j=0}^1 \sum_{i=1}^m S'_{j,i}$ .
  - it appears in set  $S'_{j,i}$ .
- b) sends  $R_{j,i}$  to  $\mathcal{C}_{j,i}$ .

3) Extracting Plaintext Intersection. Each  $\mathcal{C}_{j,i}$  takes the following steps.

- a) constructs a vector  $\vec{v}_{j,i} = [\vec{v}_{j,i,1}, \dots, \vec{v}_{j,i,m}]$ , where each vector in  $\vec{v}_{j,i}$  is initially empty.
- b) decrypts  $R_{j,i}$ 's elements as follows.  $\forall e' \in R_{j,i} :$ 
  - i) retrieves decryption key  $k_{i,l}$  and index  $l$  from  $T_{j,i}$  using  $e'$ .
  - ii) calls  $\text{PRP}^{-1}(k_{i,l}, e') \rightarrow e$  and appends  $e$  to  $l$ -th vector in  $\vec{v}_{j,i}$ .
- c) considers  $\vec{v}_{j,i}$  as the result.

Fig. 1: First Variant of Efficient Group PSI (EG-PSI<sup>(1)</sup>).

the intersection, it marks the index of the elements in the intersection and consider the elements with the same index in the unencrypted set as elements in the intersection. Similarly, each client of  $\mathcal{G}_0$  shares their encrypted set with each client of  $\mathcal{G}_1$  and follows the same steps. Ultimately, every client of  $\mathcal{G}_0$  knows the intersection between their set and all the sets belonging to the clients of  $\mathcal{G}_1$ , and the other way around. We give a detailed description of EG-PSI<sup>(1)</sup> in Figure 2.

In EG-PSI<sup>(1)</sup>,  $\mathcal{T}\mathcal{E}\mathcal{E}$  does not participate in the protocol except the OPRF evaluation and that is also a one-time activity. As a result,  $\mathcal{T}\mathcal{E}\mathcal{E}$  does not learn anything about the set

- *Parties.* Trusted execution environment  $\mathcal{T}\mathcal{E}\mathcal{E}$ , clients in group  $\mathcal{G}_0 : \{\mathcal{C}_{0,1}, \dots, \mathcal{C}_{0,m}\}$ , and clients in group  $\mathcal{G}_1 : \{\mathcal{C}_{1,1}, \dots, \mathcal{C}_{1,m}\}$ .
- *Inputs.* Sets  $S_{0,1}, \dots, S_{0,m}, S_{1,1}, \dots, S_{1,m}$ , where each  $S_{j,i}$  belongs to client  $\mathcal{C}_{j,i}$ ,  $0 \leq j \leq 1$  and  $1 \leq i \leq m$ .
- *Outputs.*  $\vec{v}_{j,i}$  to  $\mathcal{C}_{j,i}$ , where  $\vec{v}_{j,i} = \left[ [S_{j,i} \cap S_{1-j,1}], \dots, [S_{j,i} \cap S_{1-j,m}] \right]$ .

1) Setup.

$\mathcal{T}\mathcal{E}\mathcal{E}$  generates a secret key  $k$  for PRF that will be used for an OPRF evaluation.

2) Encryption of each client's set.

Each client  $\mathcal{C}_{j,i}$  and  $\mathcal{T}\mathcal{E}\mathcal{E}$  run an OPRF protocol and obtains the encrypted set  $\text{PRF}(k, S_{j,i})$ , where  $0 \leq j \leq 1$  and  $1 \leq i \leq m$ .

3) Finding Encrypted Intersections.

- a) Each client  $\mathcal{C}_{j,i}$  sends each  $\text{PRF}(k, S_{j,i})$  to every client  $\mathcal{C}_{1-j,i}$  (in the other group), where  $0 \leq j \leq 1$  and  $1 \leq i \leq m$ .
- b) Each client  $\mathcal{C}_{j,i}$ , takes the following steps.
  - i) Creates an empty set  $R_{j,i,t}$  for all  $1 \leq t \leq m$ .
  - ii) Let  $e_\ell$  represents an element of  $\text{PRF}(k, S_{j,i})$ , where  $1 \leq \ell \leq |\text{PRF}(k, S_{j,i})|$ .  
Performs  $R_{j,i,t} \leftarrow R_{j,i,t} \cup \{e_\ell\}$  if  $e_\ell$  appears both in  $\text{PRF}(k, S_{j,i})$  and  $\text{PRF}(k, S_{1-j,i})$ .

4) Extracting Plaintext Intersection.

Each  $\mathcal{C}_{j,i}$  for  $j = 0, 1$ , and  $1 \leq i \leq m$  does the following.

- a) Constructs a vector  $\vec{v}_{j,i} = [\vec{v}_{j,i,1}, \dots, \vec{v}_{j,i,m}]$ , where each vector in  $\vec{v}_{j,i}$  is initially empty.
- b) Appends the  $\ell$ -th element of  $S_{j,i}$  to  $\vec{v}_{j,i,t}$  for all  $\ell \in R_{j,i,t}$ .
- c) Returns  $\vec{v}_{j,i}$ .

Fig. 2: Second Variant of Efficient Group PSI (EG-PSI<sup>(2)</sup>).

intersection, it only learns the cardinality of each client's set.

**Theorem 2.** Let  $f_{\text{G-PSI}}$  be the functionality defined in Relation 5 (on page 4). Also, let  $\mathcal{L}$  be the leakage function presented in Definition 4 with the empty output. If OPRF is secure, then EG-PSI<sup>(2)</sup> securely realizes  $f_{\text{G-PSI}}$ , w.r.t. Definition 2.

Appendix B outlines the proof of Theorem 2. Note that both variants remain secure in the presence of semi-honest adversaries. However, if  $\mathcal{T}\mathcal{E}\mathcal{E}$  or a client is a malicious (or active) adversary, then it can affect the correctness of the result, without being detected.

#### IV. PRIVACY-PRESERVING MULTI-PARTY DEDUPLICATION (P-MPD)

##### A. Formal Definition of P-MPD

P-MPD considers the setting where there are  $n$  clients  $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$  and each  $\mathcal{C}_i \in \mathcal{C}$  has a set  $S_i$ . P-MPD enables

the clients to (efficiently) remove duplicates from their local sets such that after the deduplication the concatenation of all clients' local sets equals the union of their initial sets.

The functionality  $f_{\text{P-MPD}}$  that P-MPD computes takes a set  $S_i$  from each  $C_i$ . It returns an updated (i.e., deduplicated) set  $S'_i$  to each  $C_i$ . Formerly,  $f_{\text{P-MPD}}$  can be defined as follows:

$$f_{\text{P-MPD}}(S_1, \dots, S_m) \rightarrow (S'_1, \dots, S'_m), \quad (8)$$

where  $\sum_{i=1}^m S'_i = \bigcup_{i=1}^m S_i = S_{\cup}$  and  $S_{\cup}$  denotes the union of all initial sets and contains only unique elements.

To maintain generality, we define a leakage function, denoted as  $\mathcal{W}$ . The output of this function relies on the protocol implementing  $f_{\text{P-MPD}}$ . We assert that a protocol securely realizes  $f_{\text{P-MPD}}$  if whatever can be computed by a party in the protocol can be derived solely from its individual input and output, given the output of  $\mathcal{W}$ . Below, we formally state it.

**Definition 5** (Security of P-MPD). Let  $S = \{S_1, \dots, S_m\}$ , and  $S_i$  represent a set belonging to client  $C_i$ . Let  $\mathcal{W}$  denote a leakage function and  $f_{\text{P-MPD}}$  be the functionality defined in Relation 8. Then, a protocol  $\Psi$  securely realizes  $f_{\text{P-MPD}}$  in the presence of a static semi-honest PPT adversary  $\mathcal{A}$ , if for every  $\mathcal{A}$  in the real model, there exists a PPT Sim in the ideal model, such that for every  $C_i \in \{C_1, \dots, C_m\}$ , Relation 9 holds.

$$\{\text{Sim}_{C_i}^{\mathcal{W}}(S_i, S'_i)\}_S \stackrel{c}{\equiv} \{\text{View}_{C_i}^{\mathcal{A}, \Psi}(S)\}_S. \quad (9)$$

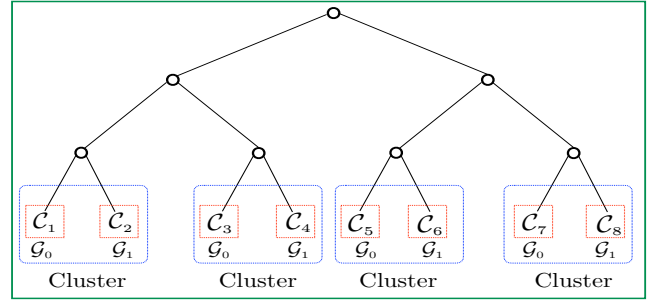
### B. Constructions of P-MPD

1) *Construction Based on EG-PSI*: We introduce a pioneering protocol called Efficient Privacy-Preserving Multi-Party Deduplication (EP-MPD) that realizes P-MPD efficiently. The idea behind EP-MPD's design involves constructing a binary tree with leaf nodes representing the clients' indices. At each level, each cluster is formed with two different groups of clients, namely  $\mathcal{G}_0$  and  $\mathcal{G}_1$ . Subsequently, EG-PSI is recursively applied to the sets of clients sharing the same cluster until the tree's root is reached. After each invocation of EG-PSI, the clients of  $\mathcal{G}_0$  update their sets by removing the intersections, returned by EG-PSI, from their local sets. These updated sets are then used as input in the subsequent EG-PSI invocation.

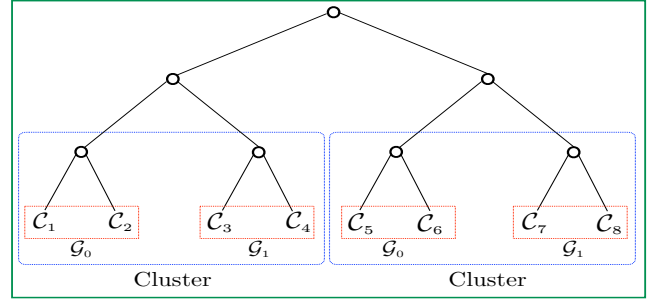
Next, we delve into further detail. We begin by sorting the clients' indices in ascending order. Next, we construct a binary tree such that the leaf nodes represent the clients' indices.

At level  $d = 1$ , commencing from the left-hand side, every two clients form a cluster. Within each cluster, the first client is assigned to group  $\mathcal{G}_0$  and the second client to group  $\mathcal{G}_1$ . Then, the clients within the same cluster engage in an instance of EG-PSI. Upon the return of the sets' intersection by EG-PSI, only the client situated on the left-hand side, specifically the one with the smaller index, modifies its local set by eliminating the elements of the intersection from its sets. This process is outlined in Figure 3a.

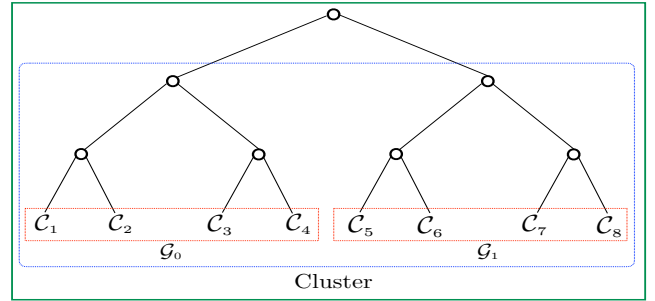
At level  $d = 2$ , starting from the left-hand side, every set of  $2^d = 2$  clients form a cluster. Within each cluster, the initial  $\frac{2^d}{2} = 2$  clients enter group  $\mathcal{G}_0$  and the remainder (i.e., the remaining 2 clients) enter group  $\mathcal{G}_1$ . Subsequently, the clients



(a) Clients joining various groups and clusters at level 1.



(b) Clients joining various groups and clusters at level 2.



(c) Depiction of how clients become part of diverse groups and clusters at level 3.

Fig. 3: The process of eight clients forming clusters and groups at various levels of a binary tree. At each level, clients belonging to groups  $\mathcal{G}_0$  and  $\mathcal{G}_1$  within the same cluster initiate EG-PSI, followed by the updating of their respective local sets.

within the same cluster engage in an instance of EG-PSI. Once again, upon the return of their sets' intersection by EG-PSI, the clients possessing the smaller indices update their local sets. This procedure is illustrated in Figure 3b. The process continues until level  $d = \log_2 m$  is reached. At this point, all  $m$  clients collectively engage in an instance of EG-PSI. Each  $\frac{2^d}{2} = \frac{m}{2}$  clients enter group  $\mathcal{G}_0$  and the remaining  $\frac{m}{2}$  clients enter group  $\mathcal{G}_1$ . These  $m$  clients jointly participate in an instance of EG-PSI. Similarly to previous steps, the clients with the smaller indices update their local sets based on the output of EG-PSI. This procedure is outlined in Figure 3c.

Figure 4 presents a detailed description of EP-MPD. For the sake of simplicity, we assume that  $m$  is a power of 2 in EP-MPD's description. However, we do not impose on the number of clients. Non-powers of 2 will result in incomplete sub-trees in the formation of clusters in the DE protocol where each group may have an unequal number of clients. This does

- *Parties.* A set of clients  $\{\mathcal{C}_1, \dots, \mathcal{C}_m\}$ .
- *Inputs.* Sets  $S_1, \dots, S_m$ , where each  $S_i$  belongs to client  $\mathcal{C}_i$ ,  $1 \leq i \leq m$ , and  $m$  is a power of 2.
- *Outputs.* Updated sets  $S'_1, \dots, S'_m$ , where each  $S'_i$  belongs to client  $\mathcal{C}_i$ , and  $\sum_{i=1}^m S'_i = \bigcup_{i=1}^m S_i$ .

The clients take the following steps  $\forall d, 1 \leq d \leq \log_2(m)$ :

- 1) *Forming Clusters.* Every distinct  $2^d$  clients enter the same cluster. Thus, for every level  $d$ , there will be  $\frac{m}{2^d}$  cluster(s). For instance, when  $d = 1$ , then there are the following clusters:  $(\mathcal{C}_1, \mathcal{C}_2), \dots, (\mathcal{C}_{m-1}, \mathcal{C}_m)$ .
- 2) *Forming Group.* In each cluster, the first  $\frac{2^d}{2}$  clients enter group 0,  $\mathcal{G}_0$ , and the rest of the clients enter group 1,  $\mathcal{G}_1$ .
- 3) *Finding Intersection.* In each cluster, the clients of  $\mathcal{G}_0$  and  $\mathcal{G}_1$  together engage an instance of EG-PSI. Each client  $\mathcal{C}_j$  receives a vector  $\vec{v}_j$  from EG-PSI, where  $\vec{v}_j$  specifies the elements that  $\mathcal{C}_j$  has in common with each client of the other group in the same cluster. It would be sufficient if only clients of  $\mathcal{G}_0$  in each cluster were aware of the intersection result.
- 4) *Removing Duplication.* For every element  $e \in \vec{v}_j$ , each  $\mathcal{C}_j$  in  $\mathcal{G}_0$  calls  $\text{Update}(S_j, \{e\}) \rightarrow S_j$ , where  $S_j \leftarrow S_j \setminus \{e\}$ .
- 5) When  $d = \log_2(m)$ , as all clients have reached the root and invoked EG-PSI at each level, they can ensure that all duplicates have been removed (due to the correctness of our EG-PSI).

Fig. 4: Efficient Privacy-Preserving Multi-Party Deduplication (EP-MPD): the construction of P-MPD.

not affect the EP-MPD protocol execution.

### C. Naive Approaches

1) *Running a Two-Party PSI Multiple Times:* One might be tempted to allow each client to invoke an existing two-party PSI with every other client. However, this approach in total requires  $\binom{m}{2} = O(m^2)$  invocations of the two-party PSI.

Whereas EP-MPD requires only  $\sum_{i=1}^{\log_2(m)} \left(\frac{m}{2^i}\right) = m - 1$  invocations of our efficient tailor-made PSI, EG-PSI. To provide concrete values, when  $m = 256$ , the naive scheme requires 32,640 invocations of a standard two-party PSI, whereas our scheme requires only 255 invocations of EG-PSI.

2) *Running Multi-Party PSI Once:* One might be tempted to execute an existing multi-party PSI protocol only once on all clients' sets, hoping to identify and subsequently remove duplicated elements. Nevertheless, this approach falls short of identifying all duplications. The limitation arises from the fact that a multi-party PSI can uncover elements common to all clients, making it incapable of detecting elements shared among only a subset of clients.

3) *Deduplicating all Sets at Once:* There is a simplified version of EP-MPD. In this variant, each client reaches an

agreement with the others on a secret key. Subsequently, each client encrypts its set elements using every key agreed upon with the other clients. Finally, all clients transmit their encrypted set elements to  $\mathcal{T}\mathcal{E}\mathcal{E}$ . Upon receiving the encrypted sets from all clients,  $\mathcal{T}\mathcal{E}\mathcal{E}$  identifies all duplicated elements. It then transmits to each client: (1) the duplicated elements found in the original elements that the respective client sent to  $\mathcal{T}\mathcal{E}\mathcal{E}$ , and (2) the index of the client with which it shares the same element. Based on  $\mathcal{T}\mathcal{E}\mathcal{E}$ 's response, the clients with smaller indices update their local sets accordingly.

However, this approach requires that  $\mathcal{T}\mathcal{E}\mathcal{E}$  maintains a local storage or memory space equal to the size of the concatenation of all encrypted sets. In contrast, EP-MPD with the underlying EG-PSI building block, allows  $\mathcal{T}\mathcal{E}\mathcal{E}$  to have a substantially smaller storage or memory space. This is achieved because, before  $\mathcal{T}\mathcal{E}\mathcal{E}$  receives all clients' encrypted sets at level  $d = \log_2(m)$ , clients apply updates to their sets multiple times. This iterative process progressively reduces the size of the clients' sets, and the reduction can be particularly significant when the number of duplications is substantial.

### D. Security Analysis

In any secure (multi-party) deduplication protocol, a participating client will learn the exact redundant elements within its set upon the protocol's completion. However, in EP-MPD, a client gains slightly more information. To be specific, upon the completion of EP-MPD, each client acquires the knowledge of the index (or identity) of the client that shares the redundant element. This additional information gained by a client during the execution of EP-MPD is formally defined as the output of a leakage function  $\mathcal{W}$  which we define below.

$\mathcal{W}$  takes as input (i) sets  $S_1, \dots, S_m$ , where each set  $S_i$  belongs to client  $\mathcal{C}_i$ , and (ii) sets  $R_1, \dots, R_m$ , where each  $R_i$  contains the redundancy (or intersection) that a set  $S_i$  has with every other set.  $\mathcal{W}$  outputs to the  $i$ -th client a set  $Q_i$  of triples, where each triple is of the form  $(r_{i,l}, i, l)$ ,  $r_{i,l} \in R_i$  is a redundancy between sets  $S_i$  and  $S_l$ .

**Definition 6.** Let  $S = \{S_1, \dots, S_m\}$  be a set of sets, where each set  $S_i$  belongs to client  $\mathcal{C}_i$ . Let  $R = \{R_1, \dots, R_m\}$  be a set of sets, where each set  $R_i$  comprises the redundancy that  $S_i$  has with every other  $j$ -th set, for all  $j$ ,  $1 \leq j \leq m$  and  $i \neq j$ . Then, leakage function  $\mathcal{W}$  is defined as:  $\mathcal{W}(S, R) \rightarrow (Q_1, \dots, Q_m)$ , where each  $Q_i$  is a set of triples of the form  $(r_{i,l}, i, l)$ ,  $r_{i,l} \in R_i$  is a redundancy between sets  $S_i$  and  $S_l$ .

**Theorem 3.** Let  $f_{\text{P-MPD}}$  be the functionality defined in Relation 8 and  $\mathcal{W}$  be the leakage function presented in Definition 6. If EG-PSI is secure (w.r.t. Definition 2), then EP-MPD (presented in Figure 4) securely realizes  $f_{\text{P-MPD}}$ , w.r.t. Definition 5.

The proof of Theorem 3 can be found in Appendix C.

### V. ENHANCED FL: APPLYING EP-MPD TO FL

With all essential building blocks in place, we are prepared to elucidate the complete functionality of the system. There are clients  $\mathcal{C}_1, \dots, \mathcal{C}_m$  having the sets  $S_1, \dots, S_m$ , where each  $S_i$  belongs to client  $\mathcal{C}_i$  respectively. Phase 1 involves each client



conducting local deduplication. Since the local deduplication is a private computation there is no privacy requirement. The data owners can apply deduplication techniques as illustrated in [2] which results in having a set  $S'_i$  free of duplicates by each client  $C_i$ . Subsequently, in Phase 2 they employ EP-MPD to eliminate duplicates across all clients. At the end of this phase, client  $C_i$  gets an updated set  $S''_i$ , such that  $\sum_{i=1}^m S''_i = \bigcup_{i=1}^m S'_i$ . Moving to Phase 3, they adopt an FL protocol (outlined in Section II-B) to train the model. Further details of this procedure are provided in Figure 5. Appendix D, proves the scheme’s correctness.

- *Parties.* A set of clients  $\{C_1, \dots, C_m\}$ .
  - *Server.* Holds the initial model  $\theta$ .
  - *Inputs.* Sets  $S_1, \dots, S_m$ , where each  $S_i$  belongs to client  $C_i$ ,  $1 \leq i \leq m$ , and  $m$  is a power of two.
  - *Outputs.* A global updated model  $\Theta$ .
- 
- 1) Local Deduplication. Each client runs a deduplication algorithm on their local dataset. At the end, client  $C_i$  receives an updated dataset  $S'_i$ .
  - 2) Global Deduplication.
    - a) All the clients participate in the EP-MPD as described in Figure 4.
    - b) Each client  $C_i$  gets updated set  $S''_i$ , such that
 
$$\sum_{i=1}^m S''_i = \bigcup_{i=1}^m S'_i.$$
  - 3) Federated learning.
    - a) The server and clients agree upon an FL protocol for training.
    - b) The server initiates the learning by sharing the initial model  $\theta$  with each client.
    - c) Each client  $C_i$  trains on their local dataset  $S''_i$  and updates  $\theta$  to  $\theta_i$ .
    - d) The clients and server aggregate the local models  $\theta_i$  trained by the clients.
    - e) The server outputs the global updated model  $\Theta$  for the next training round.
    - f) Server and clients repeat the above steps 3c–3e, until the convergence is reached.

Fig. 5: Federated learning with deduplication.

## VI. EXPERIMENTS

### A. Implementation Details

We implement the EP-MPD protocol in `python=3.10`. We use AES with a 128-bit key in CBC mode from the `cryptography` library as a PRP in EG-PSI<sup>(1)</sup>. To realize an OPRF in EG-PSI<sup>(1)</sup>, we use the OPRF algorithm proposed in [37] from the `oprflib` library. All EP-MPD experiments are conducted on a batch computing facility with a Dual Intel Xeon E5-2660 v3 @ 2.6 GHz CPU and 60 GB RAM.

For the FL experiments, we implement an FL setting for CLMs using `python=3.10` with the `transformers=4.40.1` and `torch=2.3` libraries. We implement the FedAvg [15] algorithm for FL training to analyze the effect of EP-MPD on FL’s performance. However, in practice, we recommend combining EP-MPD with privacy-preserving FL protocols for end-to-end privacy guarantees. All FL experiments are conducted on a server with an Intel(R) Xeon(R) Gold 6336Y CPU @ 2.40GHz CPU and NVIDIA RTX A6000 GPU. Appendix E contains additional experiments with large number of clients and experiments that simulate real-world network delay and bandwidth.

We run all clients sequentially on a machine for both experiments. However, to emulate a real-world setting where clients run on different devices at the same time, we first compute the individual client running times and then estimate the real-world distributed running time by assuming the clients have run simultaneously. We open-source our implementations<sup>2</sup>.

While implementing EP-MPD, we have two options for the EG-PSI module: EG-PSI<sup>(1)</sup> and EG-PSI<sup>(11)</sup>. We observe that EG-PSI<sup>(11)</sup>-based EP-MPD calls an OPRF function for every client in the leaf nodes of the EP-MPD tree, returning an encryption of the clients’ sets. As clients move to upper levels, their set elements remain unchanged, with only the cardinality potentially changing. Therefore, except at the bottom level, OPRF calls are unnecessary for encryption, allowing the same encrypted set or their subsets to be reused at other levels. Thus, we have slightly adjusted the implementation. This significantly reduces time without compromising security properties, as reflected in our experiments.

### B. EP-MPD Benchmarks

1) *Experimental Setup:* We designed three experiment settings to benchmark the performance of EP-MPD. The three settings analyze the effect of the number of clients, dataset size, and duplication percentage on the runtime of EP-MPD. All experiments were performed with both EG-PSI<sup>(1)</sup> and EG-PSI<sup>(11)</sup> as building blocks. The duplication percentage refers to the percentage of samples in a client’s dataset present in another client’s dataset. The three experiment settings are:

- a. The dataset size  $|S_i|$  is fixed at  $2^{19}$ , the duplication percentage at 30%, while the number of clients  $m$  is varied. Specifically,  $m$  is set to each value in  $\{10, 20, 30, 40, 50\}$ .
- b. The number of clients is fixed at 50 and the duplication percentage at 30%, while the dataset size  $|S_i|$  is varied. Specifically,  $|S_i|$  is set to each value in  $\{2^{10}, 2^{11}, 2^{13}, 2^{15}, 2^{17}, 2^{19}\}$ .
- c. The dataset size  $|S_i|$  is fixed at  $2^{19}$ , the number of clients  $m$  at 50, while the duplication percentage is varied. Specifically, it is set to each value in  $\{10\%, 30\%, 50\%, 70\%, 90\%\}$ .

We create pairwise duplicates between clients such that every client has a unique duplicate element with every other

<sup>2</sup><https://github.com/vdasu/deduplication>

client. All set elements are 32-bit integers. We assume that the  $\mathcal{T}\mathcal{E}\mathcal{E}$  can receive data in parallel from the clients. All other interactions that require computation between  $\mathcal{T}\mathcal{E}\mathcal{E}$  and clients occur sequentially. In the remainder of this section, by  $\text{EP-MPD}^{(1)}$  and  $\text{EP-MPD}^{(II)}$  we mean the EP-MPD that uses  $\text{EG-PSI}^{(1)}$  and  $\text{EG-PSI}^{(II)}$  respectively.

2) *Effect of Client Count:* Figure 6 depicts the effect of the number of clients on (a) EP-MPD running time, (b) the average running time of clients, and (c)  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time. Figure 6a demonstrates that the EP-MPD running time increases linearly with the number of clients for both  $\text{EP-MPD}^{(1)}$  and  $\text{EP-MPD}^{(II)}$ . The increase in the client count rises the number of PRP and OPRF evaluations as the total number of elements across all clients grows. We observe that  $\text{EP-MPD}^{(1)}$  is about 7–10 $\times$  faster than  $\text{EP-MPD}^{(II)}$ .

The average running time of the client however shows a different trend. The client running time in  $\text{EP-MPD}^{(1)}$  increases linearly with the client count and is always greater than – up to 6 $\times$  – the running time of  $\text{EP-MPD}^{(II)}$ . The running time of  $\text{EP-MPD}^{(II)}$  is constant as the number of clients increases. In  $\text{EP-MPD}^{(1)}$ , the client-side running time increases with the number of clients, because each client has to perform more PRP evaluations. However, in  $\text{EP-MPD}^{(II)}$ , the number of OPRF evaluations is equal to the number of elements in each client’s set. Additionally, in  $\text{EP-MPD}^{(1)}$ , the clients perform the PRP evaluations at every level of the tree. However, in  $\text{EP-MPD}^{(II)}$  the OPRF evaluations are only performed once when the clients are at the leaf level. The  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time in  $\text{EP-MPD}^{(II)}$  is always higher than  $\text{EP-MPD}^{(1)}$ . Both exhibit linear increases in running time as the number of clients increases. The running time in  $\text{EP-MPD}^{(II)}$  is much higher as the OPRF evaluations are relatively expensive.

3) *Effect of Dataset Size:* Figure 7 shows the effect of dataset size on EP-MPD running time, the average running time of clients, and  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time. From Figure 7a, it is evident that EP-MPD’s running time grows linearly with increasing the dataset size.  $\text{EP-MPD}^{(1)}$  is up to 8 $\times$  faster than  $\text{EP-MPD}^{(II)}$ . As the dataset size increases, the number of PRP and OPRF evaluations grows, thereby increasing the EP-MPD running time. Figure 7b shows the effect of dataset size on client running time. Unlike increasing the client count, increasing the dataset size increases the client running time of both  $\text{EP-MPD}^{(1)}$  and  $\text{EP-MPD}^{(II)}$ . We observe this difference because an increase in dataset size increases the number of OPRF evaluations while any change in the client count does not. The client running time of  $\text{EP-MPD}^{(II)}$  is up to 7 $\times$  faster than  $\text{EP-MPD}^{(1)}$ .

Figure 7c depicts the effect of dataset size on the running time of  $\mathcal{T}\mathcal{E}\mathcal{E}$ . The growth of dataset size increases the OPRF evaluations in  $\text{EP-MPD}^{(II)}$ . In  $\text{EP-MPD}^{(1)}$ ,  $\mathcal{T}\mathcal{E}\mathcal{E}$  has to identify more overlapping ciphertexts as the datasets are larger. The running time of  $\mathcal{T}\mathcal{E}\mathcal{E}$  increases linearly with the dataset size for both cases.  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time in  $\text{EP-MPD}^{(1)}$  is up to 5 $\times$  faster than  $\text{EP-MPD}^{(II)}$ .

4) *Effect of Duplication Percentage:* Figure 8 shows the effect of duplication percentage on EP-MPD running time,

the average running time of clients, and  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time. Figure 8a highlights the running time of EP-MPD as the duplication percentage changes. We note that in  $\text{EP-MPD}^{(II)}$ , OPRF evaluations incur the most overwhelming running time, dominating all other running time costs. Overall, the running time of  $\text{EP-MPD}^{(II)}$  is constant while that of  $\text{EP-MPD}^{(1)}$  linearly decreases as the duplication rate grows. In  $\text{EP-MPD}^{(II)}$ , the number of OPRF evaluations is independent of the duplication percentage, because the OPRF is always performed only when the clients are at the leaf level. However, in  $\text{EP-MPD}^{(1)}$ , the number of PRP invocations decreases as clients move up the tree from the leaf level because a higher duplication percentage removes more elements in each level of the tree. The total running time of  $\text{EP-MPD}^{(1)}$  is up to 7 $\times$  faster than  $\text{EP-MPD}^{(II)}$ .

Furthermore, the client-side running time of  $\text{EP-MPD}^{(II)}$  is constant while that of  $\text{EP-MPD}^{(1)}$  linearly decreases as the duplication percentage increases, as shown in Figure 8b. The reason is identical to the observation of the EP-MPD running time. The number of PRP evaluations reduces while the number of OPRF evaluations does not. Similar to the prior experiment settings, the client time of  $\text{EP-MPD}^{(1)}$  is higher than  $\text{EP-MPD}^{(II)}$  and is up to 7 $\times$  slower. Figure 8c shows the effect of duplication percentage on the running time of  $\mathcal{T}\mathcal{E}\mathcal{E}$ . The  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time in  $\text{EP-MPD}^{(II)}$  is constant while in  $\text{EP-MPD}^{(1)}$  linearly decreases. The  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time in  $\text{EP-MPD}^{(1)}$  is up to 2 $\times$  faster than  $\text{EP-MPD}^{(II)}$ .

The overall running time of the symmetric key-based  $\text{EP-MPD}^{(1)}$  is considerably less than the public key-based  $\text{EP-MPD}^{(II)}$ . Nevertheless, the clients in  $\text{EP-MPD}^{(1)}$  need to perform more computations than those in  $\text{EP-MPD}^{(II)}$ .

### C. Comparison with Naive Approaches

1) *Employing Fully Trusted Third Party:* One naive approach could be exploiting  $\mathcal{T}\mathcal{E}\mathcal{E}$  as a full source of trust and  $\mathcal{T}\mathcal{E}\mathcal{E}$  collects all clients’ sensitive data to build a global model. As we will explain, our approach offers both efficiency and security advantages over this naive method. Firstly, the core principle behind FL, including our scheme, is to ensure that clients retain training data within their local devices and share only model updates. Requiring clients to send plaintext data to a  $\mathcal{T}\mathcal{E}\mathcal{E}$  conflicts with this principle. While trusted execution environments provide a degree of security, they are not invulnerable—recent research has exposed their susceptibility to side-channel and other types of attacks, making them less trustworthy than previously believed.

Secondly, FL often involves many clients, each holding large datasets. Sending all data to a  $\mathcal{T}\mathcal{E}\mathcal{E}$  for model training would demand significant resources. Given the current limitations of trusted execution environments, this approach is inefficient. In contrast, FL and our system enable model training to occur on clients’ devices, relieving the computational burden on the  $\mathcal{T}\mathcal{E}\mathcal{E}$ . In our schemes, if a semi-honest adversary corrupts  $\mathcal{T}\mathcal{E}\mathcal{E}$ , then the adversary will only learn the size of the intersection that each pair of clients from a different group has (as defined in Definition 4), and nothing more, as

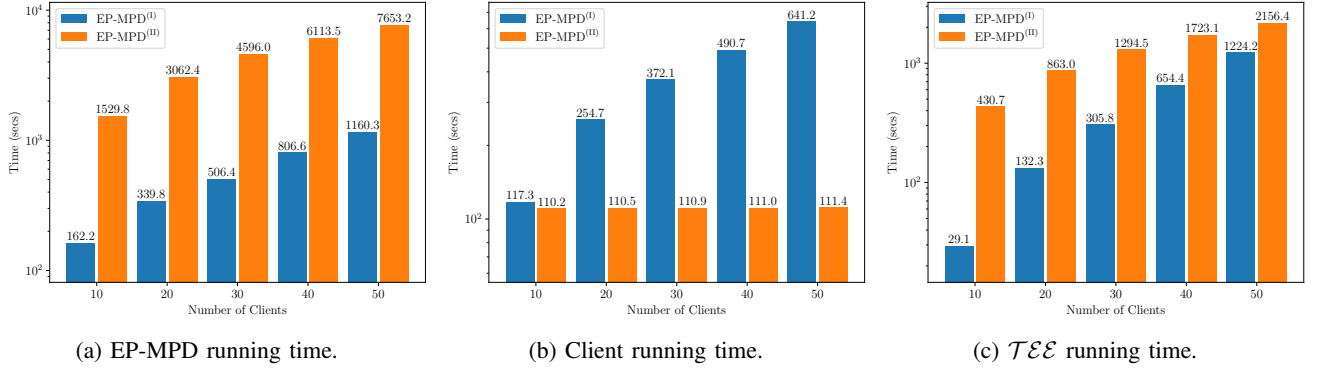


Fig. 6: Effect of client count with  $2^{19}$  dataset size and 30% duplication percentage on running time (logarithmic scale).

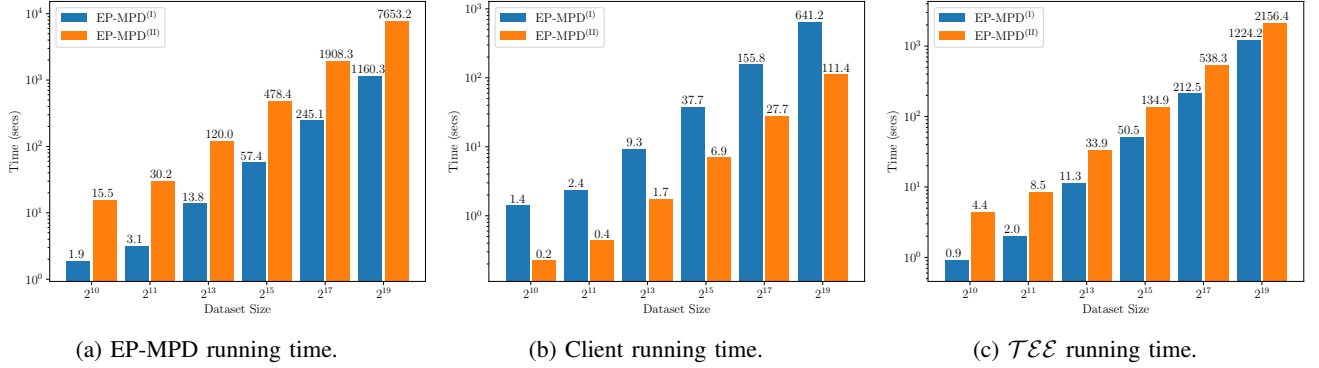


Fig. 7: Effect of dataset size with 50 clients and 30% duplication percentage on running time (logarithmic scale).

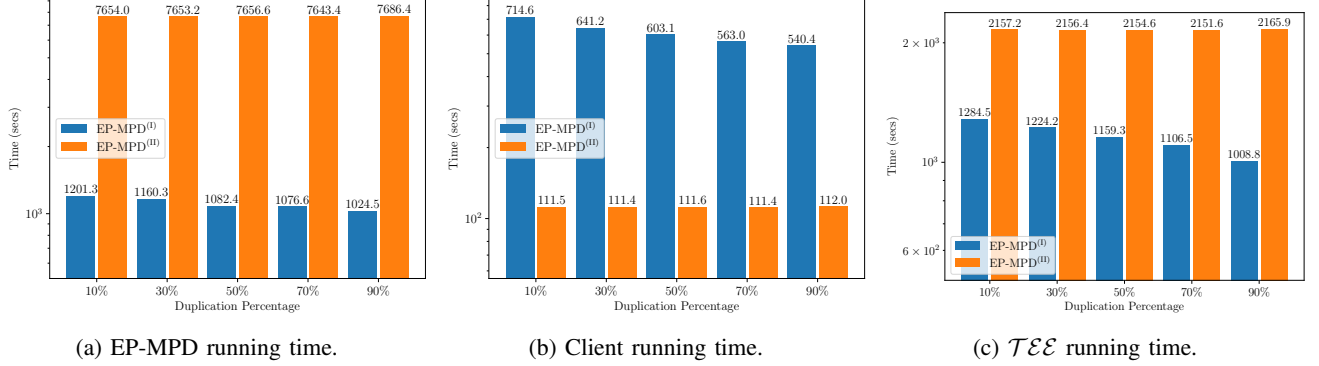


Fig. 8: Effect of duplication percentage with  $2^{19}$  dataset size and 50 clients on running time (logarithmic scale).

TABLE I: Running time (seconds) of EP-MPD and *naive* approach using two-party PSI

Client Count	EP-MPD		Naive Approach (Two-party PSI)
	EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>	
10	162.2	1529.8	616.8
20	339.8	3062.4	2537.4
30	506.4	4596.0	5890.8
40	806.6	6113.5	9884.0
50	1160.3	7653.2	15974.1

proven in Appendices A and B. Also,  $\mathcal{TE}\mathcal{E}$  never learns secret information about the parties (such as their secret keys).

2) *Comparison with Two-Party PSI*: As outlined in Section IV-C1, each client can invoke a two-party PSI protocol with every other client to remove all pairwise duplicates. However, such a naive approach would require  $\binom{n}{2}$  invocations of a two-party PSI for  $n$  clients. To highlight the practicality of EP-MPD, we implement this naive approach using the highly efficient symmetric-key based PSI proposed in [38]. Table I compares our protocols with the naive approach when the data size is  $2^{19}$  with 30% duplication. We observe that in all cases EP-MPD<sup>(I)</sup> outperforms the naive application of two-party PSI. However, the OPRF-based EP-MPD<sup>(II)</sup> is faster for settings with more than 20 clients. The speed-up offered by using

EP-MPD<sup>(1)</sup> over two-party PSI increases as the client count increases. For 10 clients we see an improvement of  $\approx 4\times$  and this increases to  $\approx 14\times$  for 50 clients.

#### D. FL with EP-MPD

1) *Setup*: We create an FL setting with varying levels of duplication to analyze the effect of EP-MPD on the model performance and training efficiency of CLM fine-tuning. We use GPT-2 Medium and GPT-2 Large CLMs from GPT-2 [39] family of decoder-only models based on the Transformer [40] architecture. GPT-2 Medium has 345 million parameters with 24 decoder blocks, 16 attention heads, and 1024 hidden sizes while GPT-2 Large has 774 million parameters with 36 decoder blocks, 20 attention heads, and 1280 hidden sizes. Both models have a sequence length of 1024 tokens.

We create 10 clients and experiment with the following 7 datasets. The Rotten Tomatoes and IMDB datasets contain movie reviews. The Sonnets and Plays datasets are a collection of works by William Shakespeare. Poetry is a collection of short poems by a variety of poets. Haiku and Short Jokes are a collection of short Japanese-style poems and Jokes from Reddit respectively. The language modeling task for each dataset is to learn how to generate text from the dataset’s distribution. For all datasets, we create a Train and Test split in an 80:20 ratio. Model perplexity is evaluated on the Test set after training. We follow a two-step procedure to create the client datasets with  $x\%$  duplication. Initially, we distributed the Train set equally among all 10 clients. Next, we sample with replacement  $x\%$  of the original Train set. This subset is then again equally distributed among all clients. We perform 5 rounds of FL training with at most 10 epochs of local training. The GPT-2 Large models were trained for fewer epochs and with less data because they were overfitting to the Train set for more epochs and data.

2) *FL model performance*: We measure the perplexity of the final model on the Test set to assess the impact of various duplication levels on model performance. Table II shows the results of perplexity (PP) and improvement rate (IR). All values in the table are rounded to two decimal places. IR refers to the percentage of improvement in perplexity after deduplication. For example, the Haiku dataset with GPT-2 Medium and 30% duplication shows an IR of 4.36% from 3.78 to 3.61 after deduplication. The deduplicated client datasets represent the baseline perplexity.

The blue text in the table highlights the largest improvement in perplexity, which is for the Plays dataset with GPT-2 Medium. We observe an improvement rate of 19.62% from 34.89 to 28.04 for 20% duplication. The Plays dataset in general has a higher perplexity than other datasets as it is a relatively challenging task with longer sequences and fewer training samples. The impact of duplicates in such scenarios is aggravated. We observe a similar trend with the Sonnets dataset with a 13.64% improvement rate for 30% duplication, which is also a difficult task with few samples and relatively long sequences.

The purple text highlights the least improvement in perplexity for the Haiku dataset with GPT-2 Medium. The 10% duplication results in a 1.09% improvement rate from 3.65 to 3.61. Similarly, the Rotten Tomatoes, Short Jokes, and IMDB datasets exhibit smaller improvement rates across different levels of duplication, while maintaining low perplexity. We observe these trends because these datasets are relatively easier to learn with a greater number of training samples. The impact of duplicates in such cases is not as pronounced.

The red text highlights two scenarios where deduplication results in insignificant changes ( $< 1\%$ ) to perplexity. The Plays dataset with GPT-2 Medium shows an improvement of 0.3% after deduplication with 10% duplication while the Sonnets dataset with GPT-2 Large worsens by 0.29% after deduplication with 10% duplication. The reason for the “noisy” results for the Plays and Sonnets dataset is that both datasets have very few training data samples. A 10% duplication percentage merely adds 3–5 duplicated data points to each client. Both datasets, however, show significant improvements for 20% and 30% duplication.

3) *FL model training efficiency*: We analyze the model training efficiency by computing the total GPU running time of all clients during training. The CPU-bound EP-MPD costs are negligible for 10 clients and at most 5,000 elements (the size of the largest dataset is 50,000 elements distributed among all clients). Table III highlights the results for GPU running time. All values in the table are rounded to two decimal places. The Time column is the total GPU running time and IR is the percentage improvement compared to the running time within the deduplicated setting.

The blue text shows the largest improvement in running time of 27.95% from 33.13 minutes to 23.87 minutes in the GPT-2 Medium experiment with the Sonnets dataset. The purple text highlights the smallest improvement in running time of 7.33% from 11.87 to 11.0 minutes in the Sonnets experiment with GPT-2 Large. We observe a relatively consistent improvement trend of about 22%, 16%, and 8% for 30%, 20%, and 10% duplication levels respectively. The GPT-2 Medium experiments show higher running times as these models were trained longer. We observe that the running time efficiency consistently improves across all duplication levels and experiments.

Tables II and III show the effectiveness of applying EP-MPD to FL of language models. The experiments show improvements of up to 19.62% in perplexity and 27.95% in GPU running time. Deduplication significantly reduces resource usage while enhancing the quality of FL models. The RTX A6000 GPU used in our experiments consumes 0.3 kWh of energy. As FL scales to larger datasets and more clients, the energy savings with EP-MPD become more significant, contributing to more sustainable and cost-effective FL implementations, benefiting the environment and the end-users.

TABLE II: Test set perplexity (PP) and improvement rate (IR) of perplexity after deduplication.

Model	Dataset	Duplication Percentage						Deduplicated
		30%		20%		10%		
		PP	IR (%)	PP	IR (%)	PP	IR (%)	PP
GPT-2 Medium	Haiku	3.78	4.36	3.7	2.37	3.65	1.09	3.61
	Rotten Tomatoes	2.4	3.62	2.36	2.0	2.35	1.67	2.31
	Short Jokes	3.96	5.34	3.89	3.79	3.83	2.31	3.74
	Poetry	6.24	14.47	6.51	18.07	5.77	7.61	5.33
	IMDB	13.17	7.1	12.57	2.67	12.49	2.05	12.23
	Sonnets	15.83	13.64	15.63	12.54	14.23	3.88	13.67
	Plays	34.32	18.3	34.89	19.62	28.12	–	28.04
GPT-2 Large	Haiku	3.26	11.29	3.25	10.83	2.98	2.78	2.89
	Rotten Tomatoes	2.65	16.79	2.61	15.29	2.53	12.81	2.21
	Short Jokes	4.11	7.84	4.03	5.86	3.94	3.64	3.79
	Sonnets	8.52	5.58	8.4	4.27	8.02	–	8.04

TABLE III: Total GPU training time (minutes) of all clients and improvement rate (IR) of time after deduplication.

Model	Dataset	Duplication Percentage						Deduplicated
		30%		20%		10%		
		Time	IR (%)	Time	IR (%)	Time	IR (%)	Time
GPT-2 Medium	Haiku	111.64	23.06	101.67	15.51	93.82	8.44	85.9
	Rotten Tomatoes	162.79	21.7	151.54	15.89	138.76	8.14	127.46
	Short Jokes	396.62	27.85	338.69	15.51	313.35	8.68	286.15
	Poetry	101.33	22.55	94.25	16.73	86.87	9.66	78.48
	IMDB	2103.93	23.99	1945.94	17.82	1772.44	9.77	1599.24
	Sonnets	33.13	27.95	28.53	16.33	26.14	8.68	23.87
	Plays	31.48	22.9	29.38	17.39	26.95	9.94	24.27
GPT-2 Large	Haiku	21.0	22.05	19.65	16.69	18.02	9.16	16.37
	Rotten Tomatoes	70.74	23.08	65.26	16.63	59.91	9.18	54.41
	Short Jokes	340.75	22.93	313.65	16.27	288.86	9.08	262.63
	Sonnets	13.91	20.92	12.89	14.66	11.87	7.33	11.0

## VII. ASYMPTOTIC COSTS ANALYSIS

In this section, we analyze the asymptotic costs of our schemes. Table IV summarizes the result.

### A. Costs of EG-PSI<sup>(l)</sup> and EG-PSI<sup>(m)</sup>

1) *Computation Cost*: Initially, we focus on the computation cost of parties in EG-PSI<sup>(l)</sup>. In step 1b, each client  $C_{j,i}$  encrypts each element of its set using the symmetric key encryption PRP, under each key it agrees with the clients in the other group. Specifically, it invokes PRP  $m \cdot |S_{j,i}|$  times, where  $m$  is the total number of clients in the other group. In step 3, each client  $C_{j,i}$  invokes PRP  $m \cdot |R_{j,i}|$  times, where  $R_{j,i}$  is the (encrypted) intersection set. Thus, the computation complexity of each  $C_{j,i}$  is  $O(m \cdot |S_{j,i}|)$ . The computation complexity of  $\mathcal{T}\mathcal{E}\mathcal{E}$  in step 2 is  $O(m \cdot |S|)$ , where  $|S|$  represents the maximum cardinality of sets. The only computation that  $\mathcal{T}\mathcal{E}\mathcal{E}$  performs is searching to find duplicated values.

Next, we analyze the computation cost of parties in EG-PSI<sup>(m)</sup>. In step 2, each  $C_{j,i}$  invokes  $|S_{j,i}|$  times OPRF to encrypt its set elements. In steps 3 and 4, each  $C_{j,i}$  performs a search on  $m+1$  encrypted sets, which imposes negligible costs

compared to the executions of OPRF. Thus, the complexity of each  $C_{j,i}$  is  $O(m \cdot |S|)$ , where  $|S|$  is the maximum cardinality of sets.  $\mathcal{T}\mathcal{E}\mathcal{E}$  invokes an instance of OPRF in step 2 for every set of elements of each client. Thus, its computation complexity is  $O(m \cdot |S|)$ , where  $|S|$  is the maximum cardinality of sets.

2) *Communication Cost*: We first focus on the communication cost of parties in EG-PSI<sup>(l)</sup>. Each  $C_{0,i}$  in step 1a transmits  $m$  keys to the clients of the other group. Thus, its communication complexity in this step is  $O(m)$ . However, the clients in  $\mathcal{G}_1$  do not transmit any message in this step. In step 1b, each  $C_{j,i}$  of each group, sends  $m \cdot |S_{j,i}|$  messages to  $\mathcal{T}\mathcal{E}\mathcal{E}$ , where the size of each message is only about 256 bits. Thus, each client's communication complexity is  $O(m \cdot |S_{j,i}|)$ .  $\mathcal{T}\mathcal{E}\mathcal{E}$  sends to each client  $R_{j,i}$  which contains the intersection the client has with the clients in the other group. Thus, the overall communication complexity of  $\mathcal{T}\mathcal{E}\mathcal{E}$  is  $O(m \cdot |R|)$ , where  $|R|$  is the maximum cardinality of sets intersection that a client may have with clients of the other group.

Now, we evaluate the communication complexity of parties in EG-PSI<sup>(m)</sup>. In step 2, each  $C_{j,i}$  performs OPRF evaluations on  $S_{j,i}$ . This procedure imposes  $O(|S_{j,i}|)$  communication cost

TABLE IV: Complexities of our solutions; namely, EG-PSI<sup>(l)</sup>, EG-PSI<sup>(m)</sup>, and EP-MPD. In the table, we have considered the maximum communication complexity of a client in EP-MPD (without considering the cases where the updates reduce the clients’ set size). In the table,  $m$  is the number of clients in each group,  $|S|$  is the maximum cardinality of sets, and  $|R|$  is the maximum cardinality of sets intersection that a client has with clients of the other group.

Protocols	Computation Cost		Communication Cost	
	Client	$\mathcal{T}\mathcal{E}\mathcal{E}$	Client	$\mathcal{T}\mathcal{E}\mathcal{E}$
EG-PSI <sup>(l)</sup>	$O(m \cdot  S_{j,i} )$	$O(m \cdot  S )$	$O(m \cdot  S_{j,i} )$	$O(m \cdot  R )$
EG-PSI <sup>(m)</sup>	$O(m \cdot  S )$	$O(m \cdot  S )$	$O(m \cdot  S_{j,i} )$	$O(m \cdot  S )$
EP-MPD	$O(m \cdot \log_2(m) \cdot  S )$	$O(m \cdot \log_2(m) \cdot  S )$	$O(m \cdot \log_2(m) \cdot  S_{j,i} )$	$O(m \cdot \log_2(m) \cdot  S )$

to the client. In step 3, the communication complexity of each  $\mathcal{C}_{j,i}$  is  $O(m \cdot |S_{j,i}|)$  because it sends its encrypted set to every client in the other group. In step 2,  $\mathcal{T}\mathcal{E}\mathcal{E}$  performs OPRF evaluations on the elements of each  $\mathcal{C}_{j,i}$  using OPRF. This procedure imposes  $O(m \cdot |S|)$  communication cost to  $\mathcal{T}\mathcal{E}\mathcal{E}$ .

### B. Costs of EP-MPD

1) *Computation Cost:* In step 3, at each level  $d$ , each  $\mathcal{C}_{j,i}$  invokes an instance of EG-PSI. When EG-PSI<sup>(l)</sup> or EG-PSI<sup>(m)</sup> is used, then the computation complexity of  $\mathcal{C}_{j,i}$  in this step is  $O(m \cdot \log_2(m) \cdot |S|)$ , where  $|S|$  is the maximum cardinality of the clients’ sets. In step 4,  $\mathcal{C}_{j,i}$  updates its set. This involves simply removing items from its local set and the associated cost is negligible compared to other costs. Thus, the overall computation complexity of  $\mathcal{C}_{j,i}$  is  $O(m \cdot \log_2(m) \cdot |S|)$ . The computation cost is dominated by the cost of executing PRP and OPRF if EG-PSI<sup>(l)</sup> and EG-PSI<sup>(m)</sup> are used respectively. The computation complexity of  $\mathcal{T}\mathcal{E}\mathcal{E}$  in either case (when EG-PSI<sup>(l)</sup> or EG-PSI<sup>(m)</sup> is used) is  $O(m \cdot \log_2(m) \cdot |S|)$ .

2) *Communication Cost:* In step 3, at each level  $d$ , each client  $\mathcal{C}_{0,i}$  transmits its updated set to an instance of EG-PSI<sup>(l)</sup> or EG-PSI<sup>(m)</sup>. Therefore, its overall communication complexity is  $O\left(m \cdot \left(|S_{0,i}| + \sum_{u=2}^d (|S_{0,i}| - \sum_{u'=1}^{u-1} |S_{\cap}^{(u')}|)\right)\right)$  where  $d = \log_2(m)$  and  $|S_{\cap}^{(u')}|$  is the size of the intersection that  $\mathcal{C}_{0,i}$  has with the clients of the other group at level  $u'$ . On the other hand, at each level  $d$ , the communication complexity of each client  $\mathcal{C}_{1,i}$  is  $O(m \cdot \log_2(m) \cdot |S_{1,i}|)$ . Note that the majority of the clients will eventually become members of group 0. Thus, their set size eventually reduces when they (i) proceed to a higher level in the tree and (ii) have intersections with other clients’ sets.

The communication complexity of  $\mathcal{T}\mathcal{E}\mathcal{E}$ , when EG-PSI<sup>(l)</sup> is used, is  $O(m \cdot \log_2(m) \cdot |R|)$ , where  $|R|$  is the maximum cardinality of sets intersection that a client may have with clients of the other group. The communication complexity of  $\mathcal{T}\mathcal{E}\mathcal{E}$ , when EG-PSI<sup>(m)</sup> is used, is  $O(m \cdot \log_2(m) \cdot |S|)$ .

## VIII. RELATED WORKS

A detailed deduplication method has been proposed in [2] that used existing techniques of finding duplicates such as suffix array [41] and MinHash [42]. However, there were no privacy requirements in this method as the data owner is responsible for local deduplication. As we consider the

problem of deduplication in a distributive setting, thus our problem statement is fundamentally different from [2].

In [43], the problem of privacy-preserving deduplication by a centralized service was discussed – there is a server that stores multiple users’ files after removing duplicates. Their method involves a second server which helps in setting up keys for each user using an OPRF. Our work is different from [43] in many ways. First, we apply OPRF for encrypting the dataset elements, which allows efficient deduplication. Second, our deduplication requirement is more granular – we want entries of users’ datasets to be deduplicated. If the two datasets have the same elements irrespective of their order, then at the end of our deduplication, one of the users will receive an empty dataset and the other will receive their dataset unchanged. On the other hand, since [43] considers a file as a whole (not every set element), so two files  $M_1$  and  $M_2$  have the same elements but in a different order will not be removed. Third, our public key-based deduplication protocol is federated as opposed to [43] which is centralized (all data is stored in a centralized server responsible for the deduplication). Fourth, [43] is a two-server setting – one server helps in deduplication and the other stores the deduplicated files. In contrast, in our protocol apart from the users, we only have one third-party.

Another work along the lines of [43] is [44] which proposed Multi-Key Revealing Encryption (MKRE) that lets a server compress the encrypted files of multiple users by removing ciphertexts related to similar plaintexts. Their construction is based on message lock encryption, therefore, it is subject to dictionary attacks [45]. Our scheme is not vulnerable to such attacks and, as a result, provides a stronger security guarantee.

The work of [46] also has pointed out the need for data processing for ML and applied PSI-type tools. However, their problem statement is to find a mismatch in the datasets belonging to two participants. In contrast, we aim to find duplicates (matching data) in sets belonging to two or more parties. Additionally, their application domain is narrow as they only considered labeled data. Their end goal is to find entries in two different sets where their features match but their corresponding labels do not. On the other hand, our protocol supports both labeled and unlabeled data.

## IX. CONCLUSIONS

It was established in [2] that machine learning benefits significantly from deduplication. However, it was not evident

how federated learning (FL) could benefit from deduplication without compromising data privacy. Our successful and efficient solution to this problem has paved the way for further applications. Our pioneering protocols, EG-PSI<sup>(1)</sup> and EG-PSI<sup>(11)</sup>, can be employed in any distributed system to identify and efficiently remove duplicates. Although this paper considers elements as duplicates only if they are exact copies, there are scenarios where elements that are not exact matches but are sufficiently similar should also be treated as duplicates [2]. Our current primitives can easily be extended to address this case. By applying fuzzy hash functions, which map closely related elements to the same digest, the deduplication of near-match elements can be reduced to the problem of deduplication of exact-match elements, allowing the use of EG-PSI<sup>(1)</sup> and EG-PSI<sup>(11)</sup> in distributed setups. While our system’s application focuses on text-based learning, it can also be extended to other areas, such as image processing and speech recognition, thereby broadening its utility and impact.

#### ACKNOWLEDGMENTS

Aydin Abadi was supported in part by EPSRC under grants EP/Y028813/1 and EP/W003325/1. Sumanta Sarkar acknowledges EPSRC grants EP/T014784/1 and EP/X036669/1.

#### REFERENCES

- [1] I. F. Ilyas and X. Chu, *Data Cleaning*. New York, NY, USA: Association for Computing Machinery, 2019.
- [2] K. Lee, D. Ippolito, A. Nystrom, C. Zhang, D. Eck, C. Callison-Burch, and N. Carlini, “Deduplicating training data makes language models better,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022.
- [3] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *J. Mach. Learn. Res.*, 2020.
- [4] N. Carlini, D. Ippolito, M. Jagielski, K. Lee, F. Tramèr, and C. Zhang, “Quantifying memorization across neural language models,” 2023. [Online]. Available: <https://arxiv.org/abs/2202.07646>
- [5] A. Suri, P. Kanani, V. J. Marathe, and D. W. Peterson, “Subject membership inference attacks in federated learning,” 2023. [Online]. Available: <https://arxiv.org/abs/2206.03317>
- [6] J. Mattern, F. Mireshghallah, Z. Jin, B. Schölkopf, M. Sachan, and T. Berg-Kirkpatrick, “Membership inference attacks against language models via neighbourhood comparison,” in *ACL*, 2023.
- [7] M. Nasr, N. Carlini, J. Hayase, M. Jagielski, A. F. Cooper, D. Ippolito, C. A. Choquette-Choo, E. Wallace, F. Tramèr, and K. Lee, “Scalable extraction of training data from (production) language models,” 2023.
- [8] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson, A. Oprea, and C. Raffel, “Extracting training data from large language models,” in *USENIX Security*, 2021.
- [9] A. Karamolegkou, J. Li, L. Zhou, and A. Sogaard, “Copyright violations and large language models,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.13771>
- [10] N. Kandpal, E. Wallace, and C. Raffel, “Deduplicating training data mitigates privacy risks in language models,” 2022. [Online]. Available: <https://arxiv.org/abs/2202.06539>
- [11] M. R. U. Rashid, V. A. Dasu, K. Gu, N. Sultana, and S. Mehnaz, “Fltrojan: Privacy leakage attacks against federated language models through selective weight tampering,” 2024. [Online]. Available: <https://arxiv.org/abs/2310.16152>
- [12] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, “On the dangers of stochastic parrots: Can language models be too big?” in *FAccT*, 2021.
- [13] E. Strubell, A. Ganesh, and A. McCallum, “Energy and policy considerations for deep learning in NLP,” in *ACL*, 2019.
- [14] D. A. Patterson, J. Gonzalez, Q. V. Le, C. Liang, L. Munguia, D. Rothchild, D. R. So, M. Texier, and J. Dean, “Carbon emissions and large neural network training,” *CoRR*, 2021.
- [15] H. B. McMahan, E. Moore, D. Ramage, and B. A. y Arcas, “Federated learning of deep networks using model averaging,” *CoRR*, 2016.
- [16] R. S. Antunes, C. A. da Costa, A. Küderle, I. A. Yari, and B. Eskofier, “Federated learning for healthcare: Systematic review and architecture proposal,” *TIST*, 2022.
- [17] J. C. Jiang, B. Kantarci, S. Oktug, and T. Soyata, “Federated learning in smart city sensing: Challenges and opportunities,” *Sensors*, 2020.
- [18] Q. Xia, W. Ye, Z. Tao, J. Wu, and Q. Li, “A survey of federated learning for edge computing: Research problems and solutions,” *High-Confidence Computing*, 2021.
- [19] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, “Federated learning for mobile keyboard prediction,” *CoRR*, 2018.
- [20] M. J. Freedman, K. Nissim, and B. Pinkas, “Efficient private matching and set intersection,” in *EUROCRYPT*, 2004.
- [21] A. Abadi, “Earn while you reveal: Private set intersection that rewards participants,” *CoRR*, 2023.
- [22] A. Abadi, C. Dong, S. J. Murdoch, and S. Terzis, “Multi-party updatable delegated private set intersection,” in *FC*, 2022.
- [23] L. Fowl, J. Geiping, S. Reich, Y. Wen, W. Czaja, M. Goldblum, and T. Goldstein, “Decepticons: Corrupted transformers breach privacy in federated learning for language models,” 2023. [Online]. Available: <https://arxiv.org/abs/2201.12675>
- [24] M. Balunović, D. I. Dimitrov, N. Jovanović, and M. Vechev, “Lamp: Extracting text from gradients with language model priors,” 2022. [Online]. Available: <https://arxiv.org/abs/2202.08827>
- [25] K. Cheng, T. Fan, Y. Jin, Y. Liu, T. Chen, D. Papadopoulos, and Q. Yang, “Secureboost: A lossless federated learning framework,” *IEEE Intelligent Systems*, 2021.
- [26] V. A. Dasu, S. Sarkar, and K. Mandal, “Prov-fl: Privacy-preserving round optimal verifiable federated learning,” ser. *AISec’22*, 2022.
- [27] R. Xu, N. Baracaldo, Y. Zhou, A. Anwar, and H. Ludwig, “Hybridalpha,” *AISec*, 2019.
- [28] J. H. Bell, K. A. Bonawitz, A. Gascón, T. Lepoint, and M. Raykova, “Secure single-server aggregation with (poly)logarithmic overhead,” ser. *CCS*, 2020.
- [29] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, “Practical secure aggregation for privacy-preserving machine learning,” in *CCS*, 2017.
- [30] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” ser. *CCS*, 2016.
- [31] O. Goldreich, *The Foundations of Cryptography - Volume 2, Basic Applications*. Cambridge University Press, 2004.
- [32] G. Chen, S. Chen, Y. Xiao, Y. Zhang, Z. Lin, and T. Lai, “Sgxpectre: Stealing intel secrets from SGX enclaves via speculative execution,” *IEEE Secur. Priv.*, 2020.
- [33] F. Dörre, J. Mechler, and J. Müller-Quade, “Practically efficient private set intersection from trusted hardware with side-channels,” in *ASIACRYPT*, J. Guo and R. Steinfeld, Eds., 2023.
- [34] R. Pass, E. Shi, and F. Tramèr, “Formal abstractions for attested execution secure processors,” in *EUROCRYPT*, 2017.
- [35] F. Zhang and H. Zhang, “Sok: A study of using hardware-assisted isolated execution environments for security,” in *HASP*, 2016.
- [36] F. Tramèr, F. Zhang, H. Lin, J. Hubaux, A. Juels, and E. Shi, “Sealed-glass proofs: Using transparent enclaves to prove and sell knowledge,” in *EuroS&P*, 2017.
- [37] J. Burns, D. Moore, K. Ray, R. Speers, and B. Vohaska, “Ec-oprff: Oblivious pseudorandom functions using elliptic curves,” *Cryptology ePrint Archive*, 2017.
- [38] S. Kamara, P. Mohassel, M. Raykova, and S. Sadeghian, “Scaling private set intersection to billion-element sets,” in *Financial Cryptography and Data Security*, N. Christin and R. Safavi-Naini, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014, pp. 195–215.
- [39] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” 2019.
- [40] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *CoRR*, vol. abs/1706.03762, 2017. [Online]. Available: <http://arxiv.org/abs/1706.03762>

- [41] U. Manber and E. W. Myers, “Suffix arrays: a new method for on-line string searches,” *SIAM J. Comput.*, 1993.
- [42] A. Z. Broder, “On the resemblance and containment of documents,” in *SEQUENCES*, 1997.
- [43] S. Keelvedhi, M. Bellare, and T. Ristenpart, “Dupless: Server-aided encryption for deduplicated storage,” in *USENIX Security*, 2013.
- [44] D. E. Lucani, L. Nielsen, C. Orlandi, E. Pagnin, and R. Vestergaard, “Secure generalized deduplication via multi-key revealing encryption,” in *SCN*, 2020.
- [45] —, “Secure generalized deduplication via multi-key revealing encryption,” Cryptology ePrint Archive, Paper 2020/799, 2020. [Online]. Available: <https://eprint.iacr.org/2020/799>
- [46] E. Blass and F. Kerschbaum, “Private collaborative data cleaning via non-equivalent psi,” in *IEEE SP*, 2023.
- [47] J. Böhrler and F. Kerschbaum, “Secure multi-party computation of differentially private median,” in *USENIX Security*, 2020.

## APPENDIX A SECURITY PROOF OF EG-PSI<sup>(1)</sup>

Next, we prove the security of EG-PSI<sup>(1)</sup>, Theorem 1.

*Proof.* We consider cases where in each case a party is corrupt.

### A. Corrupt Client

In this case, for simplicity, we focus on the view of a client  $C_{1,i}$  in  $\mathcal{G}_1$ . Because the clients of both groups behave the same way, with the main difference being that a client in  $\mathcal{G}_0$  generates a key and sends it to each client of  $\mathcal{G}_1$ . Thus, the proof asserting the protocol is secure regarding a corrupt client in  $\mathcal{G}_1$  will imply the security of the same protocol concerning a client in  $\mathcal{G}_0$ . In the real execution of the protocol, the view of a client,  $C_{1,i}$ , is defined as follows:

$$\{\text{View}_{C_{1,i}}^{A, \text{EG-PSI}^{(1)}}(S, \emptyset)\}_S = \{S_{1,i}, r_{1,i}, \vec{k}, R_{1,i}, \vec{v}_{1,i}\},$$

where  $S$  is the set containing all clients’ sets,  $r_{1,i}$  is the outcome of internal random coins of  $C_{1,i}$ ,  $\vec{k} = [k_{i,1}, \dots, k_{i,m}]$  is a vector of  $m$  keys received from the clients of  $\mathcal{G}_0$ ,  $R_{1,i}$  is the encrypted intersection  $\mathcal{C}_{1,i}$  received from  $\mathcal{T}\mathcal{E}\mathcal{E}$ , and  $\vec{v}_{j,i}$  equals  $[\vec{v}_{j,i,1}, \dots, \vec{v}_{j,i,m}]$ , each  $l$ -th vector in  $\vec{v}_{j,i}$  contains the elements that are common between  $C_{1,i}$ ’s set and the set that belongs to the  $l$ -th client in the other group.

We proceed to construct a simulator,  $\text{Sim}_{C_{1,i}}$ , in the ideal model. The simulator, receiving the client’s input  $S_{1,i}$  and output  $\vec{v}_{1,i}$ , operates as follows.

- 1) initiates an empty view. It appends to it  $S_{1,i}$  and uniformly random coins  $r'_{1,i}$ . It constructs an empty vector  $\vec{k}'$  and then chooses  $m$  keys for PRP. It appends the keys to  $\vec{k}'$  and adds  $\vec{k}'$  to the view.
- 2) constructs an empty set  $R'$ . It encrypts the elements of every  $l$ -th vector  $\vec{v}_{1,i,l} \in \vec{v}_{1,i}$  using  $l$ -th key  $k_l$  in  $\vec{k}'$  and PRP. It inserts the ciphertexts into  $R'$ .
- 3) appends  $R'$  and  $\vec{v}_{1,i}$  to the view. It outputs the view.

Now, we are prepared to discuss that the two views are computationally indistinguishable. The input set  $S_{1,i}$  has identical distribution in both models, as their elements are identical. Also, since we are in the semi-honest model, in the real execution of the protocol, the client picks its random coins  $r_{1,i}$  according to the protocol description. In the ideal model, coins  $r'_{1,i}$  have been selected uniformly at random. Therefore  $r_{1,i}$  and  $r'_{1,i}$  have identical distributions. Vectors  $\vec{k}$  and  $\vec{k}'$  have

identical distributions too because their elements have been chosen uniformly at random. In the real model,  $R_{1,i}$  contains the encrypted elements of  $\vec{v}_{1,i}$ , where the keys in  $\vec{k}$  are used for the encryption by applying the PRP. Similarly, in the ideal model, each element of  $R'$  is the encryption of an element in  $\vec{v}_{1,i}$  where a key in  $\vec{k}'$  is used for the encryption with the same PRP. Thus,  $R$  and  $R'$  have identical distributions. Also, elements of  $\vec{v}_{1,i}$  are identical in both models; therefore, they have identical distributions in the real and ideal models. We conclude that the views in the real and ideal models are indistinguishable. It is worth noting that the above proof can be easily extended to the case where  $n-1$  clients collude with each other. In this case, their *joint* views in the real and ideal models are computationally indistinguishable.

### B. Corrupt $\mathcal{T}\mathcal{E}\mathcal{E}$

In the real model, the view of  $\mathcal{T}\mathcal{E}\mathcal{E}$  is as follows:

$$\{\text{View}_{\mathcal{T}\mathcal{E}\mathcal{E}}^{A, \Gamma}(S, \emptyset)\}_S = \{(S'_{0,1}, \dots, S'_{0,m}), (S'_{1,1}, \dots, S'_{1,m})\}$$

where each  $S'_{j,i}$  contains the encrypted set elements of client  $C_{j,i}$ . Next, we construct a simulator  $\text{Sim}_{\mathcal{T}\mathcal{E}\mathcal{E}}^{\Gamma}$  that receives the output  $\vec{s}$  of leakage function  $\mathcal{L}$  (as defined in Definition 4).

- 1) initiates an empty view.
- 2) constructs  $m$  empty sets for each  $\mathcal{G}_0$  and  $\mathcal{G}_1$ . Specifically, it constructs empty sets  $\underbrace{(S''_{0,1}, \dots, S''_{0,m})}_{\mathcal{G}_0}, \underbrace{(S''_{1,1}, \dots, S''_{1,m})}_{\mathcal{G}_1}$ .
- 3) fills each set  $S''_{j,i}$  as follows, given that  $\vec{s} = [\vec{s}_{0,1}, \dots, \vec{s}_{0,m}, \vec{s}_{1,1}, \dots, \vec{s}_{1,m}]$ , where  $\vec{s}_{j,i} = \left[ |[S_{j,i} \cap S_{1-j,1}]|, \dots, |[S_{j,i} \cap S_{1-j,m}]| \right]$ , as defined in Definition 3.
  - a) for every  $y$ -th element of each  $\vec{s}_{0,i}$  (where  $1 \leq i, y \leq m$ ) picks  $\vec{s}_{0,i}[y]$  random values (from the range of PRP). It appends these values to both sets  $S''_{0,i}$  and  $S''_{1,y}$ .
  - b) pads every set  $S''_{j,i}$  with some random values (from the range of PRP) such that the size of each  $S''_{j,i}$  after padding equals the size of  $S'_{j,i}$ .
- 4) appends sets  $(S''_{0,1}, \dots, S''_{0,m}), (S''_{1,1}, \dots, S''_{1,m})$  to the view and outputs the view.

We show that the views in the real and ideal models are indistinguishable. In the real model, each elements of a set  $S'_{j,i}$  is an output of PRP. However, in the ideal model, each element of a set  $S''_{j,i}$  has been picked uniformly at random (from the range of PRP). By the security of PRP (i.e., an output of PRP is computationally indistinguishable from a random value), sets’ elements in the real and ideal model are computationally indistinguishable. The size of the intersection between each set  $S'_{0,i}$  and each set  $S'_{1,l}$  (in the real model) equals the size of the intersection between each set  $S''_{0,i}$  and each set  $S''_{1,l}$  (in the ideal model). Also, the size of each  $S'_{j,i}$  in the real model equals the size of the corresponding set  $S''_{j,i}$  in the ideal model. Hence, sets  $S'_{0,1}, \dots, S'_{0,m}, S'_{1,1}, \dots, S'_{1,m}$  in the real model are indistinguishable from sets  $S''_{0,1}, \dots, S''_{0,m}, S''_{1,1}, \dots, S''_{1,m}$  in the ideal model. Thus, the views in the real and ideal models are computationally indistinguishable.  $\square$



APPENDIX B  
SECURITY PROOF OF EG-PSI<sup>(II)</sup>

*Proof (sketch).* It is not hard to simulate parties' views.  $\mathcal{T}\mathcal{E}\mathcal{E}$ 's view in the real model can be easily simulated given that the main information it learns within the protocol execution includes its view during the invocation of the OPRF instances (but not the outputs of OPRF). If OPRF is secure, then the views of  $\mathcal{T}\mathcal{E}\mathcal{E}$  within the real and ideal model will be computationally indistinguishable in this protocol.

Each client's view (in the real execution of the protocol) mainly includes the outputs of OPRF received from its counterparts and the exchanged messages during OPRF's execution with  $\mathcal{T}\mathcal{E}\mathcal{E}$ . Due to the security of OPRF, an output of OPRF is computationally indistinguishable from a value picked uniformly at random from the output range of OPRF. Also, the messages each client receives during the execution of OPRF can be simulated, due to the security of OPRF.  $\square$

APPENDIX C  
SECURITY PROOF OF EP-MPD

In this section, we prove the security of EP-MPD, i.e., Theorem 3. For the sake of concreteness and because we have provided a detailed proof for EG-PSI<sup>(I)</sup>, we will prove Theorem 3 when EP-MPD relies on EG-PSI<sup>(I)</sup>. The proof for EG-PSI<sup>(I)</sup>-based EP-MPD will easily follow.

*Proof.* In the real execution of the protocol, the view of a client,  $\mathcal{C}_i$ , is defined as follows:  $\{\text{View}_{\mathcal{C}_i}^{\mathcal{A}, \text{EP-MPD}}(S)\}_S = \{S_i, \text{View}_{\mathcal{C}_i}, \vec{v}_i\}$ , where  $S = \{S_1, \dots, S_m\}$ ,  $S_i$  represents a set belonging to client  $\mathcal{C}_i$ ,  $\text{View}_{\mathcal{C}_i} = [\text{View}_{\mathcal{C}_i, 1}^{\mathcal{A}, \text{EG-PSI}^{(I)}}, \dots, \text{View}_{\mathcal{C}_i, z}^{\mathcal{A}, \text{EG-PSI}^{(I)}}]$  contains  $z = \log_2(m)$  views of  $\mathcal{C}_i$  such that each view  $\text{View}_{\mathcal{C}_i, d}^{\mathcal{A}, \text{EG-PSI}^{(I)}}$  corresponds to  $\mathcal{C}_i$ 's view when it invokes EG-PSI<sup>(I)</sup> at level  $d$ , where  $1 \leq d \leq z$ . Also,  $\vec{v}_i$  contains  $z$  vectors, where each vector  $\vec{v}_{i, d}$  in  $\vec{v}_i$  is the output of EG-PSI<sup>(I)</sup> that client  $\mathcal{C}_i$  receives at level  $d$ . We construct a simulator,  $\text{Sim}_{\mathcal{C}_i}$ , in the ideal model. The simulator, receives the output  $Q_i$  of leakage function  $\mathcal{W}$  (as defined in Definition 6), the client's input  $S_i$  and output  $\vec{v}_i$ .  $\text{Sim}_{\mathcal{C}_i}$  operates as follows.

- 1) initiates an empty view. For each level  $d$  ( $1 \leq d \leq z$ ) constructs the following set for each client  $\mathcal{C}_i$  in the other group, using (i) the result  $\vec{v}_{i, d}$  that  $\mathcal{C}_i$  receives at level  $d$  and (ii) the leakage function's output  $Q_i$ .
  - a) sets an empty set  $S_i$ . It finds each triple  $(x, a, b)$  in  $Q_i$ , where  $a = i$  and  $b = l$ . It adds the first element  $x$  of each triple that meets the above conditions to  $S_i$ . This ensures that  $S_i$  contains all the elements in common with  $S_i$ .
  - b) pads every  $S_i$  with random values (from the set universe) such that the size of each  $S_i$ , in the ideal model, after padding equals the size of  $S_i$  in the real model.
- 2) generates  $\frac{2^d}{2} - 1$  empty sets for the group to which  $S_i$  belongs. It appends to these empty sets values picked uniformly at random from the set universe. It ensures the size of each of these sets in the ideal model equals to that of in the real model.

3) for every level  $d$ :

- a) invokes an instance of EG-PSI<sup>(I)</sup> using the following sets: (i)  $S_i$  and sets generated in step 2, that are placed in one group and (ii) sets generated in step 1, placed in another group. As a result, it receives output  $\vec{v}_{i, d}$ .
- b) extracts  $\text{View}_{\mathcal{C}_i, d}^{\mathcal{A}, \text{EG-PSI}^{(I)}}$  after EG-PSI<sup>(I)</sup>'s execution.
- c) appends every  $\text{View}_{\mathcal{C}_i, d}^{\mathcal{A}, \text{EG-PSI}^{(I)}}$  and  $\vec{v}_{i, d}$  to its view.
- d) outputs the view.

Now, we show that the two views are indistinguishable. The input  $S_i$  in both models are identical, therefore they will have identical distribution. Due to the security of EG-PSI<sup>(I)</sup>, for every level  $d$ ,  $\text{View}_{\mathcal{C}_i, d}^{\mathcal{A}, \text{EG-PSI}^{(I)}}$  in the real and ideal models are computationally indistinguishable, as demonstrated in the proof of Theorem 1. Also, the intersection  $\vec{v}_{i, d}$  that  $\mathcal{A}$  learns at every level  $d$  in the real and ideal model are identical, thus they are indistinguishable. Hence, the views of  $\mathcal{A}$  in the real and ideal models are computationally indistinguishable.  $\square$

APPENDIX D  
PROOF OF CORRECTNESS OF THE FEDERATED LEARNING  
WITH DEDUPLICATION

In this section, we prove the correctness of the protocol presented in Figure 5, i.e., federated learning with deduplication. Informally, a protocol is correct, if it terminates and produces the expected result (i.e., a global model trained on deduplicated databases) in the presence of honest parties.

**Theorem 4.** *The “Federated learning with deduplication” scheme (presented in Figure 5) is correct if the underlying federated learning scheme, the local duplication scheme, and EP-MPD (presented in Figure 4) are correct.*

Before we prove Theorem 4, we first prove the correctness of (i) EG-PSI (i.e., EG-PSI<sup>(I)</sup> or EG-PSI<sup>(II)</sup>) which is a component we have developed that forms the foundation for EP-MPD, and (ii) EP-MPD, which is another component we have developed and upon which the federated learning with deduplication scheme relies.

**Claim 1.** EG-PSI<sup>(I)</sup> and EG-PSI<sup>(II)</sup> are correct, if PRP and OPRF are correct respectively.

*Proof:* Initially, we prove EG-PSI<sup>(I)</sup>'s correctness. Each client in  $\mathcal{G}_0$  with each client in  $\mathcal{G}_1$  agrees on an identical key and encrypts their set elements using PRP and that key. Due to PRP's correctness (i.e., its deterministic feature), if these clients have elements in common, their encryption will be identical. Thus,  $\mathcal{T}\mathcal{E}\mathcal{E}$  in Phase 2 will find the common encrypted elements of the two clients and inform them about these elements. Furthermore, due to the correctness of PRP and because each pair of clients (each belonging to a different group) that share an identical key initially store triple  $(e'_{i, l}, k_{i, l}, l)$ , given a common encrypted element, say  $e'_{i, l}$ , they can extract each related key, e.g.,  $k_{i, l}$  and decrypt the associated ciphertext, (i.e.,  $\text{PRP}^{-1}(k_{i, l}, e'_{i, l}) \rightarrow e$ ) to identify in its local database the plaintext elements common in both clients' databases. Thus, they can find the intersection of the two databases.

Now, we prove the correctness of EG-PSI<sup>(II)</sup>. Due to the correctness of OPRF (i.e., the deterministic nature of it) and because  $\mathcal{T}\mathcal{E}\mathcal{E}$  uses an identical secret key as an input of OPRF for all clients, if any client has an intersection with other clients, the output of OPRF for that element will be identical as well. This allows each client to find the encryption of identical elements in Phase 3. Since each client knows which plaintext element of its set corresponds to which ciphertext, it can identify the plaintext elements that are in common with each client of the other group. Thus, it can find the intersection of its set and its counterpart’s set.  $\square$

**Claim 2.** EP-MPD is correct if EG-PSI is correct.

*Proof:* We initially focus on the leaf node level. We know that at this level, (1) each group contains a single client and (2) EG-PSI is applied to a pair of groups. Therefore, in Phase 3, due to the correctness of EG-PSI (i.e., Claim 1), each client can always find the elements it has in common with its counterpart, i.e., the duplicated elements in these two clients’ sets. Hence, the client in  $\mathcal{G}_0$  can always delete these duplicated elements in Phase 4. At any level above the leaf node level, (1) the clients in the same group will not have any elements in common, because their intersection has already been identified and removed at the lower levels and (2) due to the correctness of EG-PSI (i.e., Claim 1) the intersection of the sets of each client in one group with each client of the other group is identified by these two clients. Thus, each client in  $\mathcal{G}_0$  can delete the duplicated element in Phase 4, at any level above the leaf node level as well. Moreover, since only a client in  $\mathcal{G}_0$  deletes the identified duplicated elements, a copy of that element still remains in the set of  $\mathcal{G}_1$ , meaning that union of  $\mathcal{G}_0$ ’s set and  $\mathcal{G}_1$ ’s set is the same after this stage.  $\square$

*Proof.* We prove Theorem 4. Due to the correctness of a local deduplication mechanism (that can simply find replicated elements in a single dataset and remove the repeated elements) each client will not have any repeated elements by the end of Phase 1. Also, due to the correctness of EP-MPD (i.e., Claim 2), by the end of Phase 2, there will be no duplicated elements across all clients’ data. EP-MPD also ensures the union of all clients’ datasets is the same after the completion of this deduplication. FL runs on the union of all clients’ datasets, hence, due to the FL’s correctness (discussed in Section II-B) the protocol will complete and all parties will always learn the final global model, trained on the deduplicated datasets.  $\square$

APPENDIX E  
ADDITIONAL EXPERIMENTS

A. Real-world Communication Setup

The experiments in Section VI-B use an implementation of the EP-MPD protocol in the ideal setting to estimate the total computation time and the computation time of the clients and  $\mathcal{T}\mathcal{E}\mathcal{E}$ . Data transfer in the ideal setting is free as the various parties directly access the data they share with

TABLE V: Running time (seconds) of EP-MPD in LAN and WAN settings

Client Count	Dataset Size	LAN		WAN	
		EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>	EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>
30	2 <sup>15</sup>	48.24	119.48	49.69	156.41
50	2 <sup>14</sup>	47.73	103.77	80.76	126.8
100	2 <sup>12</sup>	33.92	50.79	49.68	88.08
150	2 <sup>11</sup>	37.09	38.48	44.42	63.85

each other. To estimate the real-world running time, we re-implement the EP-MPD protocol using socket programming. Similar to the experiments in Section VI-B, the clients and  $\mathcal{T}\mathcal{E}\mathcal{E}$  run sequentially on the same machine, and the real-world distributed running time is estimated assuming the clients have run simultaneously. We set the network bandwidth and delay using the Linux `tc` utility and create configurations for the LAN and WAN settings. For the LAN setting, we set the network bandwidth to 1 Gbits/s and the delay to 0.02 ms round trip time (RTT). In the WAN setting, we follow the configuration in [47] where the network bandwidth is 100 Mbits/s and the delay is 100 ms (RTT). The experiments were conducted on a laptop with an Intel i7-8750H (12) @ 4.100GHz CPU, 16GB RAM, and Ubuntu 22.04.1 LTS.

Table V shows the results in the LAN and WAN settings. We fix the duplication percentage to 30% and vary the number of clients between 30 and 150. The dataset size is varied between 2<sup>11</sup> and 2<sup>15</sup>. We observe that in all cases the EP-MPD protocols efficiently deduplicate the datasets in under 3 minutes. Additionally, the overheads incurred because of the higher latency in WAN is negligible compared to the LAN.

B. Increased Client Count

TABLE VI: Running time (seconds) of EP-MPD, clients, and  $\mathcal{T}\mathcal{E}\mathcal{E}$  with large number of clients

Client Count	Total		Client		$\mathcal{T}\mathcal{E}\mathcal{E}$	
	EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>	EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>	EP-MPD <sup>(I)</sup>	EP-MPD <sup>(II)</sup>
100	718.84	3822.52	305.04	28.14	1154.38	1076.22
150	279.46	1438.31	113.25	7.13	491.97	405.34
200	531.22	1915.73	144.72	7.19	1149.73	541.25

Table VI shows the total running time, client time, and  $\mathcal{T}\mathcal{E}\mathcal{E}$  time for scenarios with a high number of clients. The dataset size is fixed to 2<sup>17</sup> for 100 clients and 2<sup>15</sup> for 150 and 200 clients. The duplication level for all experiments is set to 30%. The total number of samples in the combined datasets are  $\approx$  13 million,  $\approx$  4.9 million, and  $\approx$  6 million samples in the 100, 150, and 200 client count settings. All experiments were conducted on the same machine and experimental setup as the benchmark experiments in Section VI-B. From the table, we see that the EP-MPD protocols can efficiently scale to large number clients while also ensuring each client’s computation is low. Notably, EP-MPD<sup>(I)</sup> takes  $\approx$  530 seconds to deduplicate the dataset with  $\approx$  6 million elements held across 200 clients. With EP-MPD<sup>(II)</sup>, in the same setting, the running time of each client is just 7 seconds.

### A. Description & Requirements

The artifact is hosted online in a public GitHub repository and can be accessed at <https://github.com/vdasu/deduplication>. Additionally, the artifact is also available on Zenodo (DOI: <https://doi.org/10.5281/zenodo.14251896>).

1) *How to access*: Simply clone the GitHub or Zenodo repository on your local machine.

2) *Hardware dependencies*: The artifact comprises two separate parts. The first implements the EP-MPD protocol as a Python library. The second creates an FL setting to analyze the effect that duplicates have on GPU training time and performance. The first part of the artifact can be run on any commodity hardware. The second part requires a GPU system with at least 50GB VRAM. For example, we use a server with the NVIDIA RTX A6000 GPU.

3) *Software dependencies*: A detailed description of the software dependencies is provided in the repository in the form of Python `requirements.txt` and Anaconda `environment.yml` files.

4) *Benchmarks*: The first part of the artifact that implements the EP-MPD protocol uses synthetic data that can be generated by the library itself. The second part of the artifact that implements the FL training requires external datasets. These datasets are provided in a Google Drive link that can be found in the repository.

### B. Artifact Installation & Configuration

We provide detailed step-by-step instructions to set up the environment, install the artifact, and configure it to run the experiments in the README files in the repository.

### C. Major Claims

The major claims in the paper are:

- (C1): The EP-MPD protocol securely removes all pairwise duplicates among the datasets held by two or more clients.
- (C2): The total EP-MPD running time and  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time increase linearly, for both EP-MPD<sup>(1)</sup> and EP-MPD<sup>(m)</sup>, as the number of clients increases while the dataset size and duplication percentage are fixed. The client running time for EP-MPD<sup>(1)</sup> increases linearly while the client running time for EP-MPD<sup>(m)</sup> is constant. This claim is supported by Figure 6.
- (C3): The total EP-MPD running time,  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time, and client running increases linearly, for both EP-MPD<sup>(1)</sup> and EP-MPD<sup>(m)</sup>, as the dataset size increases while the number of clients and duplication percentage is fixed. This claim is supported by Figure 7
- (C4): The EP-MPD running time,  $\mathcal{T}\mathcal{E}\mathcal{E}$  running time, and client running remain constant for EP-MPD<sup>(m)</sup> and decrease linearly for EP-MPD<sup>(1)</sup> as the duplication percentage increases while the number of clients and dataset set is fixed. This claim is supported by Figure 8.

- (C5): The overall running time of EP-MPD<sup>(1)</sup> is less than the running time of EP-MPD two. This claim is supported by the EP-MPD running times in Figures 6, 7, and 8.

Other claims in the paper are:

- (C6): The GPU training time and perplexity of the trained models improve as the number of duplicates among the clients decreases. The duplicated setting with no duplicates offers the best performance in both, running time and model utility. This claim is supported by Tables II and III.

### D. Evaluation

Experiments E1 and E2 require a common setup. Reading the documentation in the README files in the repository and setting up the codebase requires 10 human minutes. Experiment E1 requires 1 compute minute and Experiment E2 requires 20 compute minutes for a small scale setting. More details about the feasibility of running the large-scale setting used in the paper are provided in the README. Both, the small-scale and large-scale settings, can be used to validate the claims.

1) *Experiment (E1)*: [1 compute-minute]: This experiment validates claim C1.

*[How to]* Experiment E1 simply executes the EP-MPD protocol and compares it to a naive non-private deduplication. The EP-MPD protocol has been implemented as a Python library in the artifact. The experiment first deduplicates datasets with 4 elements held by 3 clients with EP-MPD and naive non-private deduplication. Then, it compares the clients' deduplicated datasets with EP-MPD and non-private deduplication.

*[Preparation & Execution]* The README inside the EP\_MPD folder in the repository contains an example to perform EP-MPD deduplication and compares this to naive non-private deduplication. A code snippet is provided in the "Usage" section of the README file creates a setting with 3 clients and datasets with 4 elements each. The "Installation" section provides a detailed description of how to set up the environment and install the EP-MPD library.

*[Results]* The outcome of naive non-private deduplication should be the same as EP-MPD.

2) *Experiment E2*: [20 compute minutes]: This experiment validates claims C2-C4 on the effect of client count, dataset size, and duplication percentage on the total running time of the EP-MPD protocol, clients, and  $\mathcal{T}\mathcal{E}\mathcal{E}$ . Specifically, they can be used to generate Figure 6, 7, and 8. These figures are used to infer the running time trends.

*[How To]* Experiment E2 is designed to execute the EP-MPD protocol for varying configurations and generate detailed timing information during the protocol execution. The timing information is dumped to log files which are used to generate figures that validate claims C2 - C4. The generated figures need to be manually inspected to validate the claims. For example, the generated figures should have a linear increase in EP-MPD<sup>(1)</sup> client running time and constant EP-MPD<sup>(m)</sup> client running as we vary the number of clients. This is highlighted in Figure 6b.

*[Preparation & Execution]* The README inside the EP\_MPD folder in the repository provides a detailed description of how to run the experiment. The “Reproduce the paper’s results” section of the README provides steps to either generate the paper’s figures from log files or re-run the large-scale experiments in the paper. The “Small Scale Experiments” section in the README provides steps to run small-scale experiments that can be completed in under 20 minutes on commodity hardware.

*[Results]* The experiment should run without any errors and produce the log files. The generated figures should follow the same trends observed in Figure 6, 7, and 8 and validate claims C2 - C4 on the running time of EP-MPD, the clients, and  $\mathcal{T}\mathcal{E}\mathcal{E}$ .

Experiment E3 validates claim C6 and runs the FL experiments in the paper. This experiment requires a GPU with 50GB VRAM. Since the experiment deals with the federated learning of LLMs and simulates all clients on a single machine, we cannot create a small-scale version. We note that claim C6 and experiment E3 are used to motivate the need for a deduplication protocol in FL and show the effect of varying duplication percentages on FL performance. Our main contribution is the EP-MPD protocol whose claims are validated in the previous experiments.

3) *Experiment (E3)*: [10 human minutes + 4 compute-days]: This experiment fine-tunes the GPT-2 Medium and Large LMs for 10 clients using FL with various datasets and duplication percentages.

*[How to]* Experiment E3 analyzes the effect of duplicates in training data by manually implanting duplicates, with varying duplication levels, in datasets commonly used to fine-tune LLMs. The experiment fine-tunes LLMs with four duplication levels for each dataset and LLM. Namely, “Deduplicated”, “10%”, “20%”, and “30%”. The “Deduplicated” setting corresponds to the case when a deduplication protocol like EP-MPD is used to remove duplicates before training. The experiment generates log files that store the GPU training time and model perplexity. Each combination of duplication level, dataset, and LLM needs to run separately. The paper considers 44 different combinations, and experiment E3 encompasses all of them.

*[Preparation & Execution]* The artifact uses an Anaconda environment to manage dependencies and install the required packages. The README inside the FL folder in the repository provides a detailed description of how to run the experiment under the “Installation Instructions” and “Usage” sections.

*[Results]* The generated log files can be used to create Table II and III. The generated log files should highlight the same trends observed in Table II and III. The GPU training time should increase as the number of duplicates increases. Similarly, the model perplexity should worsen (increase) as the number of duplicates increases.

### E. Customization

All experiments can be customized by choosing a different set of parameters from the ones used in the paper. The README files in the repository provide a detailed description of how to run the experiments with different parameters.