

Delegated Private Matching for Compute

Dimitris Mouris
University of Delaware
jimouris@udel.edu

Daniel Masny
Meta Inc.
daniel.masny@rub.de

Ni Trieu
Arizona State University
nitrieu@asu.edu

Shubho Sengupta
Meta Inc.
ssengupta@meta.com

Prasad Buddhavarapu
Meta Inc.
bprasad@meta.com

Benjamin Case
Meta Inc.
bmcase@meta.com

ABSTRACT

Private matching for compute (PMC) establishes a match between two datasets owned by mutually distrusted parties (C and P) and allows the parties to input more data for the matched records for arbitrary downstream secure computation without rerunning the private matching component. The state-of-the-art PMC protocols only support two parties and assume that both parties can participate in computationally intensive secure computation. We observe that such operational overhead limits the adoption of these protocols to solely powerful entities as small data owners or devices with minimal computing power will not be able to participate.

We introduce two protocols to *delegate* PMC from party P to untrusted cloud servers, called *delegates*, allowing multiple smaller P parties to provide inputs containing identifiers and associated values. Our *Delegated Private Matching for Compute* protocols, called DPMC and D_s PMC, establish a join between the datasets of party C and multiple delegators P based on multiple identifiers and compute secret shares of associated values for the identifiers that the parties have in common. We introduce a rerandomizable encrypted oblivious pseudorandom function (OPRF) primitive, called EO, which allows two parties to encrypt, mask, and shuffle their data. Note that EO may be of independent interest. Our D_s PMC protocol limits the leakages of DPMC by combining our EO scheme and secure three-party shuffling. Finally, our implementation demonstrates the efficiency of our constructions by outperforming related works by approximately $10\times$ for the total protocol execution and by at least $20\times$ for the computation on the delegators.

KEYWORDS

Oblivious pseudorandom function, private identity matching, private record linkage, secure multiparty computation

1 INTRODUCTION

Cloud computing has become a prominent solution for storage and analytics since it enables clients to outsource their data and not have to worry about scalability, data availability, and most importantly maintaining their own infrastructure. Gathering data from multiple input providers and computing statistics over all of their data enables a plethora of useful applications such as gathering


real-time location data and notifying users of possible exposure to highly infectious diseases [25, 52]. In certain applications, linking client data to proprietary information owned by multiple larger entities may unlock unique insights that are otherwise not possible. Users should not have to trust that the cloud servers will not store their sensitive personal data for use for purposes other than it's intended. The problem of computing meaningful analytics across multiple input parties while preserving user privacy from cloud server providers becomes significantly more challenging.

Secure multi-party computation (MPC) offers prominent cryptographic solutions for jointly computing on private data from multiple input providers [27, 36, 54]. Although general-purpose MPC frameworks [6, 21, 29, 32, 53] enable running arbitrary computations over private data (e.g., medical data analytics [26]), they generally incur significant performance overheads compared to solutions that are tailored to one application (e.g., machine learning [33] and statistics [7, 15, 20, 40, 41]). Similarly, specialized private set intersection (PSI) protocols [12–14, 24, 34, 45, 46, 48–50] introduce significantly more efficient solutions than generic MPC but focus solely on private matching and disregard associated metadata.

A few recent protocols inspired by [37] that are based on the hardness of Decisional Diffie–Hellman (DDH) have attempted to securely link private records and allow general-purpose secure computation on the common data. More specifically, private matching for compute (PMC) [8, 10, 42] from Meta, private set intersection (PSI) [4] from Apple, and private join and compute (PJC) [31, 35] from Google enable computing intersections and unions between two parties while protecting the privacy of the underlying users. After the private linkage is computed, these protocols enable downstream secure computation based on the matched records. Unfortunately, prior works only focus on two parties and require both of them to actively participate in the private matching, which restricts the adoption of these protocols to solely powerful entities as non-crypto-savvy data owners or devices with minimal computing power will not be able to engage in secure computation protocols.

We motivate our work by focusing on the example of *ad attribution*, a crucial business application for tracking the effectiveness of online advertising campaigns and increasing revenue generated by ads. An ad conversion refers to the situation where a user interacts with an online ad for a particular product on the ad publisher's website (which we call party C) and then goes on to make a purchase on the advertiser's website (which we call party P_t for $t \in \{1, \dots, T\}$). Of course, a large ad publisher may host ads from multiple advertisers, each of which would like to know how their ad campaigns are performing and which purchases can be attributed to their online ads. However, the data required to compute these statistics are

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.



Proceedings on Privacy Enhancing Technologies 2024(2), 49–72
© 2024 Copyright held by the owner/author(s).
<https://doi.org/10.56553/popets-2024-0040>

split across multiple parties: the ad publisher knows the users who have seen a particular ad, and each advertiser knows who made a purchase and what they spent.

Additionally, in the case of multiple advertisers, the secure ad attribution protocol needs to be repeated for each advertiser, which can be both inefficient and not allow for elaborate statistics. To address this, we propose a *delegated setting* where all advertisers securely delegate the computation to a delegate party, which computes the ad attribution securely with the ad publisher. This approach reduces the computational burden from individual advertisers and also allows for computing cross-advertiser connections, which leads to more advanced applications such as personalization. Users that appear in multiple advertiser datasets are combined into the same row of the final join instead of having multiple joins between the ad publisher and each advertiser. Now the ad publisher and the delegate party can run a secure downstream computation over all advertisers’ data. In a real-world instantiation, the ad publisher (e.g., Google, Meta) can collaborate with non-profits (e.g., ISRG) or other companies (e.g., Cloudflare, Mozilla) and allow advertisers to delegate their computation.

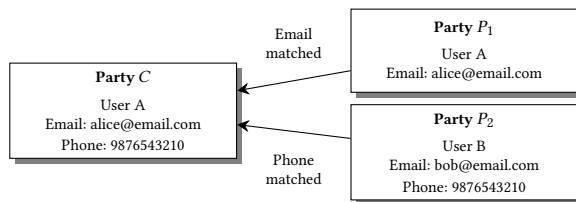


Figure 1: Many to many connections. P_1 and P_2 have different views about User A and B, whereas Party C has listed them as one user. Notably, P_1 and P_2 might also be the same party.

To improve match rates, our work considers multi-key matching [8], which involves matching users across multiple identifiers such as email, phone number, and other personal information. By using multiple keys, we can achieve a more accurate match and increase the likelihood of finding a connection between the user who saw the ad and the user who made the purchase. As shown in Fig. 1, User A may be indexed by both an email address and a phone number in party C, whereas, in other data owners, the combination of identifiers may differ (e.g., name, email). In Section 5, we demonstrate the applicability of our work in several domains such as private ad attribution, privacy-preserving analytics, and private machine learning. Note that our approach is not limited to these specific applications and can be extended to other use cases as well.

We ask the following motivating question:

How can we design delegated private record linkage protocols that allow multiple entities to outsource their records and perform secure computation on the associated metadata of the matched records?

1.1 Previous Works

We now discuss relevant works for private record linkage protocols that allow for computing over associated data. A comparison of our protocols with related works can be found in Table 1.

Private Matching for Compute (PMC) introduced DDH-based constructions for private matching that compute a union of two

datasets held by mutually distrusting parties C and P without revealing which items belong to the intersection [8, 10]. After the matching phase, both parties can input associated data for each row in the union and engage in a downstream secure computation. The core idea of PMC is to have each party first hash their records and then exponentiate them to random secret scalars. After exchanging the hashed and exponentiated records, each party exponentiates the other party’s records to their secret scalar and they both arrive at the same random identifiers.

Multi-key PMC [8] assumes that each record may have multiple identifiers (as shown in Fig. 1) and leverages a ranked deterministic join logic that collapses many-to-many connections and achieves a one-to-one mapping. The idea is that each identifier has a predefined weight and the matching is first performed on all the records based on the first identifier (as the single-key PMC) before continuing to the other identifiers. PMC leaks the intersection size to the two parties and in case of multiple keys per row, they learn the bipartite graph of matches up to an isomorphism. Additionally, PMC supports computing on the match between two datasets and requires both parties to actively participate in both the matching and the downstream secure computation, which significantly limits the adoption in real-world applications. Contrary to PMC, our work allows matching between any number of parties and shifts the cost away from the parties by delegating the computation to a server.

Private secret-shared set intersection (PS^3I) [10] is a natural extension of PMC that allows the two parties to input associated data to the matching protocol. Instead of learning a mapping to original inputs, the two parties only learn additive secret shares of those records which they can feed into any general-purpose MPC framework. PS^3I is realized using Paillier additive homomorphic encryption (HE) scheme [44] and incurs significant performance overheads. Additionally, PS^3I only works between two parties and requires both parties to be online for the whole protocol execution.

Private Join and Compute (PJC) [31, 35] computes the intersection between two datasets and aggregates the associated data for all the rows in the intersection using additive HE. Contrary to our work that computes secret shares for all the associated data, PJC only allows computing a sum of the data in the intersection. Furthermore, as with all the previous related works, PJC only supports two parties whereas our protocols scale to multiple parties.

Mohassel et al. [39] utilize cuckoo hash tables and perform SQL-like queries over two secret shared databases in the honest majority three-party setting. Both the input and output tables are secret shared between the computing parties. Because the cuckoo hash tables do not support duplicates, [39] has leakages in the presence of non-unique identifiers. Additionally, the join protocols focus on two parties and in order to compute joins between multiple parties the protocol has to be iterated multiple times. Each party’s database has to be joined with the output of the previous join or they can be combined in a binary-tree-like structure. Contrary, our delegated protocols are designed to support multiple delegators and do not have to be repeated for each input party.

Circuit-PSI relies on oblivious Pseudorandom Functions (OPRF) for computing PSI between two parties and then computing a function over the common data [12, 46, 47, 49]. Both parties learn secret shares of 1 or 0, representing record presence in the intersection. These shares are used to input associated data as both

Table 1: Comparisons with related works in terms of functionality, number of parties, threat model, and multi-key support.

Protocol	Functionality	Input Parties		Computing Parties		Delegated	Multi-key
		No.	Model	No.	Model		
PMC [10]	Join (Union)	2	Semi-honest	Same as Input Parties		✗	✗
MK-PMC [8]	Join (Union)	2	Semi-honest	Same as Input Parties		✗	✓
PS ³ I [10]	Join (Union) & Secret Share Associated Data	2	Semi-honest	Same as Input Parties		✗	✗
PJC [31]	PSI-Sum (Sum Associated Data in Intersection)	2	Semi-honest	Same as Input Parties		✗	✗
Circuit-PSI [12, 46, 47, 49]	PSI & Secret Share Associated Data	2	Semi-honest	Same as Input Parties		✗	✗
Catalic [25]	PSI-CA (Count Items in Intersection)	2	Semi-honest	≥ 2	Semi-honest	✓	✗
DB Joins [39]	Join and Select Statements	2	Semi-honest	3	Semi-honest	✓	✗
DPMC	Join (Left-Join) & Secret Share Associated Data	$T + 1$	Semi-honest	2	Semi-honest	✓	✓
D_sPMC	Join (Left-Join) & Secret Share Associated Data	$T + 1$	Semi-honest	3	Semi-honest	✓	✓

parties actively participate in the protocol. These works focus on the two-party setting and it is not clear how to extend them to the delegated setting where multiple parties outsource their data for matching along with encrypted associated data to a helper party. One could imagine that the input parties compute hash tables but then they would need to privately combine these tables which goes beyond what has been studied in the literature. On the other hand, Catalic [25] uses OPRFs between two parties but allows one party to delegate its computation to a powerful server. All the aforementioned works allow two-party matching which is solely based on a single key, while our work supports matching based on datasets of multiple parties where each can have multiple keys (e.g., name, email). Finally, our protocols enable multiple input parties to delegate their computation and then go offline instead of requiring them to participate in expensive protocols.

Miao et al. [38] introduced a shuffled distributed OPRF (DOPRF) for computing PSI between two parties and the shuffling is performed in the clear by one of the parties. Similarly, the authors of [5] propose a DDH-based PRF combined with ElGamal encryption that allows for shuffling by one of the parties. We introduce an encrypted oblivious pseudorandom function (OPRF) primitive, called EO, which allows two parties to encrypt, mask, and shuffle their data. While [5, 38] seem to be similar to our EO primitive, we emphasize that they are quite different. First of all, our EO primitive performs shuffling under MPC for security contrary to their shuffled protocol. Additionally, although EO could be instantiated with a combination of ElGamal and DH-OPRF, *EO is an abstraction* that can fit many possible instantiations (e.g., from codes, lattices, isogenies). Finally, in this work, we are in the delegated setting to allow multiple parties to outsource the private matching and go offline instead of solely focusing on the two-party setting.

1.2 Our Contributions

Delegated Protocols. We propose a new family of *Delegated Private Matching for Compute* protocols, called DPMC and D_sPMC, that build upon PMC [8, 10] and lift the burden of engaging in secure computation from parties with less computational power. Our protocols rely on a powerful server (which we call party *C*) and on a *delegate* node (which we refer to as party *D*) to perform private record linkage between the records of *C* and input parties, which we call *delegators* or parties *P*. Contrary to previous works that focus on linking data only between two parties, our work enables linkage

between *C* and multiple delegators (P_1 to P_T) and aims to make the computation in the delegators lightweight to foster wide-scale adoption. Parties *C* and *D* engage in a two-party computation to compute a private left join of party’s *C* and all the delegators’ data.¹ *C*’s input is a multi-key dataset where each row contains multiple identifiers (i.e., keys) that can be matched. The delegators’ inputs are multi-key datasets with associated data, which comprise both identifiers and metadata that can be in any form (e.g., numbers, strings). *C* learns a mapping from its users to the left join but does not learn which of its users have been matched. For each row in the left join, both *C* and *D* receive secret shares that correspond to the delegators’ associated data if that row maps to one of the delegators’ identifiers or a share of NULL (i.e., zero), otherwise.

Left Join. Our motivation for performing left join compared to a union or an intersection is that party *C* learns a mapping from all their users into the join, which allows them to input additional associated metadata without re-executing the matching protocol and without learning which users matched or not. These data can either be labels (in the clear) that could be used to filter the secret shared values (e.g., in a GROUP BY fashion), or they can be additional secret shares for the downstream MPC computation. After the matching process and the secret shares have been established, parties *C* and *D* only need to know the relative order of their shares, which can then be used for any downstream secure computation such as privacy-preserving analytics and machine learning. Our goal in this work is to create efficient protocols that can be realized in real-world applications for private left join and allow the delegators to outsource the computation to delegates.

Threat Model. We assume *semi-honest* security, which we prove in the Appendices. Party *C* follows the protocol specification but tries to exfiltrate information about the delegators’ data. Similarly, the delegate *D* tries to exfiltrate information about all parties’ data. Finally, delegators are semi-honest and outsource their real data.

Multi-key Datasets. To increase matching rates, we adopt multi-key datasets [8] and consider matching between records on more than one identifier (i.e., key) that inherently generates many-to-many connections. We use a ranking-based technique to collapse multiple connections into one-to-many (*C*-to-*P*) connections. Although we do not claim this contribution, it is an important feature

¹Our core protocol computes the left join between party’s *C* data and all delegators’ data. We show in Appendix E how to modify it to compute the inner join.

that increases match rates and complicates the protocols; note that it is not straightforward to add this to related works. We operate over unique keys (e.g., email) to avoid inference attacks [43].

Rerandomizable Encrypted OPRF (EO). We introduce a concept called EO that captures the tasks of encrypting identifiers, shuffling ciphertexts, homomorphically evaluating a PRF on the ciphertexts, and decrypting a homomorphically evaluated ciphertext to the PRF output with the identifier as input. More specifically, one party can evaluate the OPRF both on clear data (e.g., x) and encrypted data (e.g., $\text{Enc}(x)$), and then delegate the matching to another party that can decrypt the evaluated OPRF of the encrypted data and get $\text{Dec}(\text{PRF}(\text{Enc}(x))) = \text{PRF}(x)$ (keys are omitted here, see Section 3.3). Further, EO allows shuffling encrypted inputs such that the third party cannot correlate PRF outputs and the initially received ciphertexts. Notice that EO is more powerful than an OPRF since it allows encrypting inputs for the PRF and sends the ciphertexts to a third party (i.e., delegate the evaluation) whereas an OPRF asks that the input provider directly interacts with the PRF evaluator. This allows us to reduce the leakage to the third party (i.e., Party D). Furthermore, the PRF evaluation can be distributed between Party C , who owns the key and homomorphically evaluates the PRF, and Party D , who decrypts the homomorphically evaluated ciphertext and obtains the PRF output.

In Section 3.3, we define the EO primitive and we provide an instantiation based on DDH and ElGamal in Appendix B. Note that EO is an abstraction and can fit various instantiations (possibly from codes, lattices, etc.). Finally, the EO construction might be of independent interest and can facilitate other protocols as well.

D_S PMC Protocol. We use our EO primitive to extend DPMC to D_S PMC, a protocol that uses two delegates (party D and a shuffler S). D_S PMC performs an honest majority shuffling protocol between parties C , D , and S and achieves stronger security guarantees in the case of a corruption of Party D and multiple delegators.

Applications. We envision multiple applications that may leverage our delegated setup of merging multiple private datasets and securely computing analytics on metadata. A healthcare provider holding patient records may gain critical insights such as calculating the risk of a health condition by merging with data stored on individual smart devices or other healthcare providers, without needing to access identifiable user data. An ad publisher holding user-provided information may be able to measure advertising efficacy and offer personalized ads by merging with data held by millions of businesses while still preserving user privacy.

Our contributions are summarized as follows:

- We introduce a novel DPMC protocol for securely computing left join between multiple distrusting parties.
- Design of a new rerandomizable encrypted OPRF (EO) primitive that enables encrypting inputs, shuffling ciphertexts, homomorphically evaluating a PRF on ciphertexts, and decrypting ciphertexts to PRF outputs. EO is of independent interest.
- We combine EO and secure three-party shuffling to extend DPMC to D_S PMC, a protocol that reduces DPMC's leakage.
- We detail applications in online advertising such as privacy-preserving ad attribution, analytics, and personalization.

2 PRELIMINARIES

2.1 Notation

We denote the computational security parameter by κ . We use $[m]$ to refer to the set $\{1, \dots, m\}$. We denote the concatenation and exclusive OR (XOR) of two-bit strings x and y by $x \parallel y$ and $x \oplus y$, respectively. We use $r \xleftarrow{\mathbb{R}}$ to refer to a randomly chosen element r from set \mathbb{R} . We use ppt to denote probabilistic polynomial time. We use $\{\}$ for unordered and $()$ for ordered sets.

2.2 Definitions

DEFINITION 1 (MULTI-KEY KEY-VALUE STORE). A multi-key key-value store KV is a set of key sets c_i , i.e., $KV := \{c_i\}_{i \in [m]}$. Each set c_i contains m_i keys, i.e., $c_i := \{c_{i,j}\}_{j \in [m_i]}$. When the key set is ordered, we denote it with $c_i := (c_{i,j})_{j \in [m_i]}$. Further, KV might contain m values v_i for $i \in [m]$, one associated with each key set. In this case, we denote the key sets as $KV := \{p_i, v_i\}_{i \in [m]} = \{\{p_{i,j}\}_{j \in [m_i]}, v_i\}_{i \in [m]}$. Note, we use p_i instead of c_i when the set includes associated data v_i . Furthermore, each key $c_{i,j}$ in a set KV is unique, i.e., there does not exist an $(i', j') \neq (i, j)$ s.t. $c_{i',j'} = c_{i,j}$.

Informally, a multi-key left join with associate data between KV_C and KV_P results in a set of values with as many rows as the set on the left (i.e., KV_C in our case) and the associated values of KV_P for the rows that matched and zero, otherwise.

DEFINITION 2 (MULTI-KEY LEFT JOIN WITH ASSOCIATED DATA BETWEEN TWO KEY-VALUE STORES). Let $KV_C := \{c_i\}_{i \in [m_C]}$ be a multi-key set of party C that contains ordered key sets, i.e., $c_i := (c_{i,j})_{j \in [m_{C,i}]}$. Let $KV_P := \{p_i, v_i\}_{i \in [m_P]}$ be a multi-key key-value set of P that contains both key sets, i.e., $p_i := \{p_{i,j}\}_{j \in [m_{P,i}]}$ and associated values v_i . The left join between KV_C and KV_P is defined by $KV_C \bowtie KV_P := (\hat{v}_i)_{i \in [m_C]}$, where $\hat{v}_i := v_{i'}$ s.t. j_i is the smallest element in $[m_{C,i}]$ for which there exists an $i' \in [m_P]$ and $j_{i'} \in [m_{P,i'}]$ with $c_{i,j_i} = p_{i',j_{i'}}$. If there does not exist such an j_i , i' and $j_{i'}$, we define $\hat{v}_i := 0$.

With multiple delegators, we extend Def. 2 as follows.

DEFINITION 3 (MULTI-KEY LEFT JOIN WITH ASSOCIATED DATA BETWEEN $T + 1$ KEY-VALUE STORES). Let for all $t \in [T]$, $KV_t := \{p_{t,i}, v_{t,i}\}_{i \in [m_t]}$ be a multi-key key-value set of party P_t that contains both key sets, i.e., $p_{t,i} := \{p_{t,i,j}\}_{j \in [m_{t,i}]}$, and values, i.e., $v_{t,i}$. Also, let $KV_C := \{c_i\}_{i \in [m_C]}$ be a multi-key set of party C that contains only ordered key sets, i.e., $c_i := \{c_{i,j}\}_{j \in [m_{C,i}]}$. The left join between KV_C and $\{KV_t\}_{t \in [T]}$ is defined as:

$$KV_C \bowtie \{KV_1, \dots, KV_T\} := (\pi_i(\hat{v}_{i,1}, \dots, \hat{v}_{i,T}))_{i \in [m_C]},$$

where for each $t \in [T]$ and $i \in [m_C]$, $\hat{v}_{i,t}$ is defined as follows. Let for each $j \in [m_{C,i}]$, $S_{i,j,t} := \{i' \in [m_t] \mid \exists j' \in [m_{t,i'}] \text{ s.t. } c_{i,j} = p_{t,i',j'}\}$. If $\bigcup_j S_{i,j,t} \neq \emptyset$, we define $j_{i,t} := \min\{j \in [m_{C,i}] \text{ s.t. } S_{i,j,t} \neq \emptyset\}$, i' is defined as the unique $i' \in S_{i,j_{i,t},t}$ and $\hat{v}_{i,t} := v_{t,i'}$. If $\bigcup_j S_{i,j,t} = \emptyset$, we define $\hat{v}_{i,t} := 0$. Finally, the values $\hat{v}_{i,1}, \dots, \hat{v}_{i,T}$ are permuted by a random permutation π_i for each row $i \in [m_C]$.

This definition ensures that value $\hat{v}_{i,t}$ is associated with delegator t such that each row in the join corresponds to T values, one for each delegator. There might be multiple possible matching rows for each delegator with one of the identifiers in c_i . In that case, we include the row that matches with $c_{i,j}$ with the smallest j in the

join. Since each identifier is unique in each set, there is only one identifier that matches with $c_{i,j}$.

We adjust Def. 3 in case values cannot be assigned to specific delegators anymore. Therefore a row might contain multiple values of the same delegator while other delegators might not be represented with a value. Changing the definition of the join allows us to reduce the overall leakage for the D_s PMC protocol.

DEFINITION 4 (MULTI-KEY LEFT JOIN WITH ASSOCIATED DATA AND MINIMAL LEAKAGE BETWEEN $T + 1$ KEY-VALUE STORES). Let for all $t \in [T]$, $KV_t := \{p_{t,i}, v_{t,i}\}_{i \in [m_t]}$ be a multi-key key-value store of party P_t that contains both key sets, i.e., $p_{t,i} := \{p_{t,i,j}\}_{j \in [m_{t,i}]}$, and values, i.e., $v_{t,i}$. Also, let $KV_C := \{c_i\}_{i \in [m_C]}$ be a multi-key set of party C that contains only ordered key sets, i.e., $c_i := \{c_{i,j}\}_{j \in [m_{C,i}]}$. The left join between KV_C and $\{KV_t\}_{t \in [T]}$ is defined as:

$$KV_C \bowtie \{KV_1, \dots, KV_T\} := (\pi_i(\hat{v}_{i,1}, \dots, \hat{v}_{i,T}))_{i \in [m_C]},$$

where for each $i \in [m_C]$, $\hat{v}_{i,t}$ is defined as follows. Let for each $j \in [m_{C,i}]$, $S_{i,j} := \{(t', i') \in ([T], [m_{t'}]) \mid \exists j' \in [m_{t',i'}] \text{ s.t. } c_{i,j} = p_{t',i',j'}\}$. Further, we define the set of indices that have not been included in the join yet as $S_{i,j,<t} := S_{i,j} \setminus \{(t', i') \in ([T], [m_{t'}]) \mid \exists t'' < t \text{ s.t. } \hat{v}_{i,t''} = v_{t',i'}\}$. If $\bigcup_j S_{i,j,<t} \neq \emptyset$, we define $j_{i,t} := \min(j \in [m_{C,i}] \text{ s.t. } S_{i,j,<t} \neq \emptyset)$, (t', i') is defined as a random $(t', i') \stackrel{R}{\leftarrow} S_{i,j_{i,t},<t}$ and $\hat{v}_{i,t} := v_{t',i'}$. If $\bigcup_j S_{i,j,<t} = \emptyset$, we define $\hat{v}_{i,t} := 0$. Finally, the values $\hat{v}_{i,1}, \dots, \hat{v}_{i,T}$ are permuted by a random permutation π_i for each row $i \in [m_C]$.

Def. 4 defines the join as follows. It defines value $\hat{v}_{i,t}$ by matching $c_{i,j}$ for the smallest j with a match that has not yet been included in the join and takes the value of a random row i' of a random delegator t' that matches with $c_{i,j}$. If there is no such a match left, $\hat{v}_{i,t}$ is defined as 0.

We provide Algs. 1 and 2 for Defs. 3 and 4 in Appendix A.

DEFINITION 5 (KEY ENCAPSULATION MECHANISM (KEM)). A key encapsulation with security parameter κ is a triplet of algorithms $(\text{KEM.KG}, \text{KEM.Enc}, \text{KEM.Dec})$ with the following syntax.

- $\text{KEM.KG}(1^\kappa)$: On input 1^κ output a key pair $(\text{KEM.pk}, \text{KEM.sk})$.
- $\text{KEM.Enc}(\text{KEM.pk})$: On input KEM.pk , KEM.Enc outputs an encapsulation KEM.cp and key KEM.k .
- $\text{KEM.Dec}(\text{KEM.sk}, \text{KEM.cp})$: On input $(\text{KEM.sk}, \text{KEM.cp})$, KEM.Dec outputs a key KEM.k .

For correctness, we ask that

$$\Pr[\text{KEM.Dec}(\text{KEM.sk}, \text{KEM.cp}) = \text{KEM.k}] \geq 1 - \text{negl},$$

where the probability is taken over $(\text{KEM.pk}, \text{KEM.sk}) \leftarrow \text{KEM.KG}(1^\kappa)$ and $(\text{KEM.cp}, \text{KEM.k}) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$.

We need simulatable KEMs, which is true for commonly used KEMs. A KEM is simulatable if there exists a ppt algorithm KEM.Sim with $\text{KEM.cp} \leftarrow \text{KEM.Sim}(\text{KEM.sk}, \text{KEM.k})$, where KEM.cp is computationally indistinguishable from $\text{KEM.cp}' \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ under the constraint that $\text{KEM.k} = \text{KEM.Dec}(\text{KEM.sk}, \text{KEM.cp}')$. Further, we need standard key indistinguishability.

DEFINITION 6 (KEY INDISTINGUISHABILITY). We call a KEM key indistinguishable if for any ppt algorithm \mathcal{A} ,

$$\begin{aligned} & \left| \Pr[\mathcal{A}(\text{KEM.pk}, \text{KEM.cp}, \text{KEM.k}) = 1] - \right. \\ & \left. \Pr[\mathcal{A}(\text{KEM.pk}, \text{KEM.cp}, u) = 1] \right| \leq \text{negl}, \end{aligned}$$

where $(\text{KEM.pk}, \text{KEM.sk}) \leftarrow \text{KEM.KG}(1^\kappa)$, $(\text{KEM.cp}, \text{KEM.k}) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ and $u \leftarrow \{0, 1\}^*$.

DEFINITION 7 (SECRET SHARING). We call two values $sh_1, sh_2 \in \{0, 1\}^*$ a two-out-of-two XOR secret sharing of a secret value a if $sh_1 \oplus sh_2 = a$ and for $i \in \{0, 1\}$ sh_i is uniform and independent of a .

Secret sharing schemes allow a dealer to distribute shares of her data to multiple parties so that each share does not reveal anything about the original data [2]. In MPC, each party creates secret shares of their data and shares them with the other parties. Then, each party computes a function of the shares and combines them to reconstruct the final output. MPC utilizes secret sharing to compute arbitrary arithmetic functions as arithmetic circuits [2, 19, 32, 51]. In this work, we utilize binary (XOR) secret sharing as in Def. 7, but our protocols can also support arithmetic shares. To compute arbitrary functions as arithmetic circuits, XOR shares can be converted to arithmetic as in [16, 33].

DEFINITION 8 (IND-CPA SECURITY). We call an encryption scheme indistinguishable under chosen plaintext attacks (IND-CPA secure) if for any ppt algorithm \mathcal{A} ,

$$\left| \Pr[\mathcal{A}(\text{pk}, \text{ct}_0) = 1] - \Pr[\mathcal{A}(\text{pk}, \text{ct}_1) = 1] \right| \leq \text{negl},$$

where $(\text{pk}, \text{sk}) \leftarrow \text{PKE.KG}(1^\kappa)$, $(x_0, x_1) \stackrel{R}{\leftarrow} \mathcal{X}$, $\forall i \in \{0, 1\} : \text{ct}_i \leftarrow \text{PKE.Enc}(\text{pk}, x_i)$. In case of a symmetric key encryption, we replace \mathcal{A} 's access to pk with access to an encryption oracle for key sk . We also replace $(\text{PKE.KG}, \text{PKE.Enc}, \text{PKE.Dec})$ with $(\text{SKE.KG}, \text{SKE.Enc}, \text{SKE.Dec})$.

We include additional definitions such as the DDH assumption, pseudorandom generator, random oracle, and symmetric and public key encryption in Appendix A.

2.3 Ideal Functionality for Delegated PMC

We present the ideal functionality $\mathcal{F}_{\text{DPMC}}$ for Delegated PMC in Fig. 2. In the ideal world, $\mathcal{F}_{\text{DPMC}}$ is composed of a functionality for join $\mathcal{F}_{\text{JOIN}}$ and a functionality for compute \mathcal{F}_{CMP} . $\mathcal{F}_{\text{JOIN}}$ gets input from party C a multi-key set KV_C and from parties P_1 to P_T multi-key key-value sets KV_1, \dots, KV_T with associated values v_1, \dots, v_T and computes a left join \mathcal{J} with associated data as described in Def. 3 (or alternatively Def. 4). That is, for each record c_i and for each party $t \in [T]$, \mathcal{J} holds $\hat{v}_{i,t}$ which represents either the associated metadata (if there was a match) or a zero (if no match was found for c_i) as $(\hat{v}_{i,1}, \dots, \hat{v}_{i,T})_{i \in [m_C]}$. Next, $\mathcal{F}_{\text{JOIN}}$ samples secret shares SH_C and SH_D such that $\mathcal{J} = \text{SH}_C \oplus \text{SH}_D$ and sends SH_C to party C and SH_D to D . Later, C and D can query \mathcal{F}_{CMP} with their secret shares and \mathcal{F}_{CMP} first reconstructs $\mathcal{J} := \text{SH}_C \oplus \text{SH}_D$ and then computes $y := f(\mathcal{J})$. Our ideal functionality $\mathcal{F}_{\text{DPMC}}$ is composed of the functionality of join $\mathcal{F}_{\text{JOIN}}$ and compute \mathcal{F}_{CMP} .

Parties P_1 to P_T do not get any output from $\mathcal{F}_{\text{JOIN}}$ or \mathcal{F}_{CMP} , whereas C and D learn secret shares of the associated data of matched values from $\mathcal{F}_{\text{JOIN}}$. Finally, C learns the output y from \mathcal{F}_{CMP} which depends on function f . Even in the ideal world, if f returns all the associated values without performing any computation (e.g., aggregation), C does not learn which value corresponds to which user, or even which of the users in KV_C have been matched.

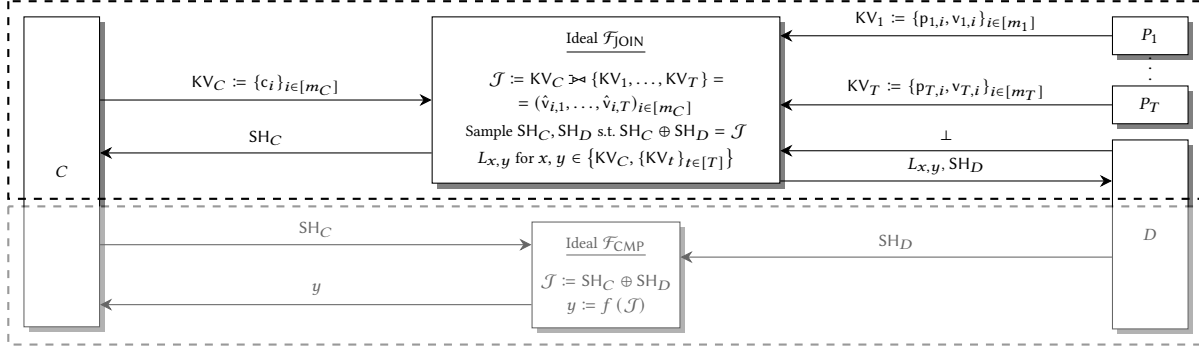


Figure 2: Ideal functionality $\mathcal{F}_{\text{DPMC}}$ of private join for compute is composed by $\mathcal{F}_{\text{JOIN}}$ and \mathcal{F}_{CMP} . Parties C and P_1, \dots, P_T provide inputs. $\mathcal{F}_{\text{JOIN}}$ computes a left join with associated data $\mathcal{J} := \text{KV}_C \bowtie \{\text{KV}_1, \dots, \text{KV}_T\} = (\hat{v}_{i,1}, \dots, \hat{v}_{i,T})_{i \in [m_C]}$ as described in Def. 3 (or alternatively Def. 4). Later on, C and D can query \mathcal{F}_{CMP} with their secret shares (SH_C and SH_D) and \mathcal{F}_{CMP} will reconstruct $\mathcal{J} := \text{SH}_C \oplus \text{SH}_D$ and compute $y := f(\mathcal{J})$ and send it to party C . Party D gets leakage $L_{x,y}$, where x and y are the sets of any party in $\{\text{KV}_C, \text{KV}_1, \dots, \text{KV}_T\}$.

Party C	Party P ₁	Party P ₂
williamfulmore@example.net (345) 678-9012	carljohanson44@example.com - \$225	cindymeiners@example.com (901) 234-5678 \$20
cindymeiners@example.net (901) 234-5678	annelopez82@example.net (234) 813-1908 \$250	...
lanastasiades@example.com (123) 456-7890	...	cpaynter@example.com (567) 605-936 \$200
...	sebastian@example.com (214) 654-1312 \$100	
carljohanson44@example.com (890) 123-4567		

(a) Input parties multi-key key-value sets containing emails, phone numbers, and dollar amounts.

XOR AD Share	Party C	Join	Party D	AD	XOR AD Share
10110110	williamfulmore@example.net (345) 678-9012	←	←	\$0	10110110
00000111	cindymeiners@example.net (901) 234-5678	←	←	\$20	00010011
00010011	lanastasiades@example.com (123) 456-7890	←	←	\$0	00010011
10001010	carljohanson44@example.com (890) 123-4567	←	←	\$225	01101011

(b) Multi-key left join (result of $\mathcal{F}_{\text{JOIN}}$). All records of C are matched. The records of P_1 and P_2 that do not match with C do not appear in the left join.

Figure 3: Multi-key left-join overview. Parties C and D compute a left-join of the multi-key sets of C , P_1 , and P_2 and the XOR secret shares of the associated data (AD) of the delegators (P_1 and P_2). In (a), we show an example with three parties (C , P_1 , and P_2). P_1 and P_2 have associated data (shown as \$ amounts; note that they might have more associated data). In (b), parties C and D have performed the left-join and ended up with secret shares of the associated data of the matched records (shown in blue). For readability, we show the associated data (AD) in (b) to indicate the value of the XOR shares, note that this remains secret.

Finally, D learns a leakage function $L_{x,y}$, where x and y each represent any of KV_C or $\text{KV}_1, \dots, \text{KV}_T$. We extend our $\mathcal{F}_{\text{DPMC}}$ functionality to $\mathcal{F}_{\text{D_sPMC}}$ that limits the aforementioned leakage between sets KV_C and KV_P , where $\text{KV}_P := \{\text{KV}_1, \dots, \text{KV}_T\}$. We provide a formal definition of the leakage function when introducing the different protocols. In the case of a single identifier per row, the leakage corresponds to the cardinality of the intersection.

Note that the MPC functions computed in the \mathcal{F}_{CMP} phase should be carefully considered for privacy reasons. Without making any assumptions about the inputs, using differential privacy seems to be the only option to protect individual users being signaled out. We discuss realistic applications in Section 5.

3 LEFT JOIN DELEGATED PMC PROTOCOLS

3.1 Overview

Our goal is to join records that represent the same entities across datasets that are held by different parties without revealing any information about the individual records. We focus on performing a left join between the datasets of multiple parties and computing secret shares of the associated data of the matched records so they can be fed into downstream general-purpose MPC.

We realize $\mathcal{F}_{\text{DPMC}}$ with two novel protocols that compute left join with associated data. We introduce a *delegate* party D that enables multiple *delegators* (parties P_1, \dots, P_T) to securely delegate their data and go offline, similarly to $\mathcal{F}_{\text{JOIN}}$ in Fig. 2. Parties C and D engage in our proposed delegated protocols to privately link

Setup: All parties agree on a g be a generator of a cyclic group \mathbb{G} with order q where DDH is hard and hash functions $H_{\mathbb{G}}(\cdot) : \{0, 1\}^* \rightarrow \mathbb{G}$, $H(\cdot) : \{0, 1\}^* \rightarrow \{0, 1\}^{|\nu_{t,i}|}$. All parties P_t have access to the public key pk_D of party D , and party D has secret key sk_D .

<p>① Key-Generation (Party C)</p> <p>1: $(\text{KEM.pk}, \text{KEM.sk}) \leftarrow \text{KEM.KG}(1^k)$</p> <p>Send to P_t and D: KEM.pk</p> <hr/> <p>② Identity Match (Party P_t)</p> <p>Input: $\text{KV}_t = \{(p_{t,i}, \nu_{t,i})\}_{i \in [m_t]}$ for data set size m_t.</p> <p>Messages: KEM.pk</p> <p>1: $a_t \xleftarrow{R} \mathbb{Z}_q$, $sk_t \leftarrow \text{SKE.KG}(1^k)$ ▷ Random scalar a_t and secret key sk_t.</p> <p>2: $\text{cta}_t := \text{PKE.Enc}(pk_D, sk_t)$, $\text{ctb}_t := \text{SKE.Enc}(sk_t, a_t)$</p> <p>3: For $i \in [m_t]$: ▷ For each row in KV_t.</p> <p>4: $(\text{KEM.cp}_{t,i}, \text{KEM.k}_{t,i}) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ ▷ New encs.</p> <p>5: $\text{ha}_{t,i} := H_{\mathbb{G}}(p_{t,i})^{a_t}$ ▷ Hash and exponentiate to a_t.</p> <p>6: $\text{sh}_{C,t,i} := \text{KEM.k}_{t,i}$ ▷ Share of $\nu_{t,i}$ for party C.</p> <p>7: $\text{sh}_{D,t,i} := \nu_{t,i} \oplus \text{sh}_{C,t,i}$ ▷ Share of $\nu_{t,i}$ for party D.</p> <p>8: $\text{ctc}_{t,i} := \text{SKE.Enc}(sk_t, (\text{KEM.cp}_{t,i}, \text{sh}_{D,t,i}))$</p> <p>Send to C: $\text{cta}_t, \text{ctb}_t$ and $\{(\text{ha}_{t,i}, \text{ctc}_{t,i})\}_{i \in [m_t]}$</p> <hr/> <p>③ Identity Match (Party C)</p> <p>Input: $\text{KV}_C = \{c_i\}_{i \in [m_C]}$ for data set size m_C.</p> <p>Messages: $\{\text{cta}_t, \text{ctb}_t, \{(\text{ha}_{t,i}, \text{ctc}_{t,i})\}_{i \in [m_t]}\}_{t \in [T]}$</p> <p>1: $a_C \xleftarrow{R} \mathbb{Z}_q$ ▷ Random scalar a_C.</p> <p>2: For $t \in [T]$, $i \in [m_t]$: ▷ For each P_t and each row.</p> <p>3: $\text{hca}_{t,i} := (\text{ha}_{t,i})^{a_C}$ ▷ Hash and exponentiate P_t's data to a_C.</p> <p>4: Pick random permutation π, $\hat{i} := \pi(t)$.</p> <p>5: For $i \in [m_C]$: ▷ For each row in KV_C.</p> <p>6: $\text{hc}_i := H_{\mathbb{G}}(c_i)^{a_C}$ ▷ Hash and exponentiate own data to a_C.</p> <p>Send to D: $\{\text{hc}_i\}_{i \in [m_C]}$ and $\{\text{cta}_t, \text{ctb}_t, \{(\text{hca}_{t,i}, \text{ctc}_{t,i})\}_{i \in [m_t]}\}_{t \in [T]}$</p>	<p>④ Identity Match and Recover Shares (Party D)</p> <p>Messages: $\text{KEM.pk}, \{\text{hc}_i\}_{i \in [m_C]}, \{\text{cta}_t, \text{ctb}_t, \{(\text{hca}_{t,i}, \text{ctc}_{t,i})\}_{i \in [m_t]}\}_{t \in [T]}$</p> <p>1: For $\hat{i} \in [T]$: ▷ For each delegator party $P_{\hat{i}}$.</p> <p>2: $sk_{\hat{i}} := \text{PKE.Dec}(sk_D, \text{cta}_{\hat{i}})$</p> <p>3: $a_{\hat{i}} := \text{SKE.Dec}(sk_{\hat{i}}, \text{ctb}_{\hat{i}})$</p> <p>4: For $i \in [m_{\hat{i}}]$: ▷ For each row in $\text{KV}_{\hat{i}}$.</p> <p>5: $(\text{KEM.cp}_{\hat{i},i}, \text{sh}_{D,\hat{i},i}) := \text{SKE.Dec}(sk_{\hat{i}}, \text{ctc}_{\hat{i},i})$</p> <p>6: $\text{hc}_{\hat{i},i} := \text{hca}_{\hat{i},i}^{1/a_{\hat{i}}}$ ▷ Remove $a_{\hat{i}}$.</p> <p>7: Join $\mathcal{J} := (\text{hc}_i)_{i \in [m_C]} \bowtie \llbracket (\text{hc}_{t,i})_{t \in [T], i \in [m_t]} \rrbracket$ ▷ Details in Alg. 1.</p> <p>8: For each row i and delegator \hat{i} in \mathcal{J}: ▷ For each row in the join.</p> <p>9: If record matched:</p> <p>10: $\widehat{\text{KEM.cp}}_{\hat{i},i} := \text{KEM.cp}_{\hat{i},i'}$ ▷ Use encs. from $P_{\hat{i}}$.</p> <p>11: $\widehat{\text{sh}}_{D,\hat{i},i} := \text{sh}_{D,\hat{i},i'}$ ▷ Use share of $\nu_{\hat{i},i}$ generated by $P_{\hat{i}}$.</p> <p>12: Else: ▷ no match found</p> <p>13: $(\widehat{\text{KEM.cp}}_{\hat{i},i}, \text{KEM.k}_{\hat{i},i}) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ ▷ New encs.</p> <p>14: $\widehat{\text{sh}}_{D,\hat{i},i} := \text{KEM.k}_{\hat{i},i}$ ▷ Use share of 0.</p> <p>15: Pick m_C random permutations $\{\pi_i\}_{i \in [m_C]}$.</p> <p>16: $\mathcal{J}_D := (\pi_i(\widehat{\text{sh}}_{D,\hat{i},i})_{\hat{i} \in [T]})_{i \in [m_C]}$ ▷ D's permuted XOR shares.</p> <p>Send to C: $\{\pi_i(\widehat{\text{KEM.cp}}_{\hat{i},i})_{\hat{i} \in [T]}\}_{i \in [m_C]}$</p> <hr/> <p>⑤ Recover Shares (Party C)</p> <p>Input: KEM.sk</p> <p>Messages: $\{\widehat{\text{KEM.cp}}_{\hat{i},i}\}_{i \in [m_C], \hat{i} \in [T]}$</p> <p>1: For $i \in [m_C]$, $\hat{i} \in [T]$: ▷ For each row and each delegator.</p> <p>2: $\widehat{\text{sh}}_{C,\hat{i},i} := \text{KEM.Dec}(\text{KEM.sk}, \widehat{\text{KEM.cp}}_{\hat{i},i})$ ▷ Get $\text{KEM.k}_{\hat{i},i}$.</p> <p>3: $\mathcal{J}_C := (\widehat{\text{sh}}_{C,\hat{i},i})_{i \in [m_C], \hat{i} \in [T]}$ ▷ Aligned with $(c_i)_{i \in [m_C]}$</p>
---	--

Figure 4: Single-key DPMC. Party C and the delegators P_1 to P_T compute the left-join of their records with the help of D . Parties C and D receive \mathcal{J}_C and \mathcal{J}_D , respectively. These sets contain XOR secret shares for each row in the join. For each delegator P_t , if a row is in the intersection, the parties hold XOR shares of the delegator's associated data, otherwise, XOR shares zero. Party C additionally learns a mapping from its users into the join a but does not learn which of its users have been matched.

C 's and all the delegators' records and compute secret shares of the associated data for the matched records. Both C and D learn T secret shares for each row in the left join that corresponds to either the delegators' associated data (i.e., ν_t) if that row maps to a record in KV_t or a secret share of zero (if that row is only in KV_C). C also receives a mapping from its users into the join but does not learn which of its users have been matched.

Having the secret shares as the protocol output allows parties C and D to realize \mathcal{F}_{CMP} and jointly compute a function f on the secret shared associated data. An intuition of our delegated protocols is shown in Fig. 3. In Fig. 3 (a), we show the multi-key datasets of party C and two delegators P_1 and P_2 . In Fig. 3 (b), we show the matching performed on both e-mail addresses and phone numbers (Def. 3), as well as the generated XOR shares. Interestingly, our protocols are compatible with both XOR and arithmetic secret shares. To keep things simple, we use XOR shares exclusively. Note that in Fig. 3 (b) we show the AD for readability – Party D does not learn the associated data (only the secret shares of them).

3.2 Delegated PMC (DPMC)

For simplicity, we start with a strawman DPMC protocol that does not operate over multi-key databases (e.g., has only email addresses). Our first variant for left join between $T + 1$ databases is shown

in Fig. 4 and consists of three stages: “key generation”, “identify match”, and “recover shares”. Both parties C and D learn a left join size (m_C) set of XOR shares (\mathcal{J}_C and \mathcal{J}_D , respectively) for each row in the join that corresponds to the delegators' associated data if that row maps one of parties' P_1 to P_T records or a secret share of zero (if that row is only present in KV_C). Additionally, C receives a mapping from its users into \mathcal{J}_C but does not learn which of its users are in the intersection. The two parties can use the secret shares \mathcal{J}_C and \mathcal{J}_D for any general-purpose MPC computation.

Intuitively, the protocol works as follows. Each party P_t hashes its identifiers $p_{t,i}$ (for each row i) using $H_{\mathbb{G}}$ and masks them with a random a_t . The associated values $\nu_{t,i}$ are secret shared where the share for C (i.e., $\text{sh}_{C,t,i}$) is the key of a KEM. Each party encrypts the shares for Party D , the mask a_t and the key encapsulation towards party D using pk_D , and sends it to C . It also sends the masked hashes of the identifiers (i.e., $H_{\mathbb{G}}(p_{t,i})^{a_t}$) to C . Note that this does not leak any information to C since $H_{\mathbb{G}}(p_{t,i})^{a_t}$ could be seen as a PRF evaluation and is therefore pseudorandom based on DDH.

Party C permutes the messages and uses a random $a_C \xleftarrow{R} \mathbb{Z}_q$ to compute $\text{hca}_{t,i} := H_{\mathbb{G}}(p_{t,i})^{a_t \cdot a_C}$. a_C can be seen as a PRF key. It forwards the permuted messages including $\text{hca}_{t,i}$ and sends the PRF evaluation of its own identifiers, i.e., $\text{hc}_i := H_{\mathbb{G}}(c_i)^{a_C}$ to D .

Party D decrypts all the ciphertexts and un.masks $H_{\mathbb{G}}(p_{t,i})^{a_t \cdot ac}$ to $H_{\mathbb{G}}(p_{t,i})^{ac}$ using a_t . It then matches the results with the hc_i 's sent by Party C . If there is a match, it just forwards the key encapsulation $\text{KEM.cp}_{\hat{i},i'}$ from delegator P_t and uses the decrypted share $\text{sh}_{D,t,i}$ as its own share. Otherwise, it generates a new encapsulation and uses the generated key $\text{KEM.k}_{i,\hat{i}}$ as its own share. In this step, we do not leak C 's share and therefore value $v_{t,i}$ to Party D , due to the key indistinguishability of the key encapsulation.

In the final step, C uses the secret key of the key encapsulation received by D to recover its own shares. Observe that for the unmatched records, we get secret shares of zero as $\mathcal{J}_C \oplus \mathcal{J}_D = \widehat{\text{sh}}_{C,i,\hat{i}} \oplus \widehat{\text{sh}}_{D,i,\hat{i}} = \text{KEM.k}_{i,\hat{i}} \oplus \text{KEM.k}_{i,\hat{i}} = 0$, while for the matched records we get secret shares of the delegators associated data as $\mathcal{J}_C \oplus \mathcal{J}_D = \widehat{\text{sh}}_{C,i,\hat{i}} \oplus \widehat{\text{sh}}_{D,i,\hat{i}} = \text{KEM.k}_{i,\hat{i}} \oplus \text{sh}_{D,i,\hat{i}} = \text{KEM.k}_{i,\hat{i}} \oplus v_{t,i} \oplus \text{sh}_{C,i,\hat{i}} = v_{t,i}$. Party D cannot distinguish shares of $v_{t,i}$ from shares of 0 since the encapsulations generated by parties P_1 to P_T have the same distribution as the ones generated by D .

Leakage. We define DPMC's leakage in Def. 9, where D learns the sizes of the intersection between each two parties. For instance, for parties C , P_1 , and P_2 , party D will learn $|\text{KV}_C \cap \text{KV}_1|$, $|\text{KV}_C \cap \text{KV}_2|$, and $|\text{KV}_1 \cap \text{KV}_2|$ but without knowing which party is P_1 and which is P_2 due to the permutation performed by C . With multiple keys, D will also learn a graph of matches as defined by $L_{x,y}$ next. We give a formal security theorem (Theorem 10) and prove it in Appendix D.1.

DEFINITION 9 (DPMC LEAKAGE). Given KV_C and $\text{KV}_1, \dots, \text{KV}_T$, the leakage $L_{x,y}$ of the ideal functionality in Fig. 2 for the DPMC protocol in Fig. 4 is defined as follows. Define $\text{KV}_{u,C}$ by replacing $c_{i,j} \in \text{KV}_C$ with $u_{i,j} \xleftarrow{R} \{0, 1\}^K$. Define $\text{KV}_{u,t}$ by replacing $p_{t,i,j} \in \text{KV}_t$ with $u'_{i',j'}$ if there exist t', i', j' with $p_{t,i,j} = c_{i',j'}$ or an already replaced $p_{t',i',j'}$ with $p_{t,i,j} = p_{t',i',j'}$, otherwise replace it with $u'_{t,i,j} \xleftarrow{R} \{0, 1\}^K$. $L_{x,y} := \{(C, \text{KV}_{u,C}), \pi(t, \text{KV}_{u,t})_{t \in [T]}\}$.

THEOREM 10. Let the secret key encryption and the PKE scheme be IND-CPA secure, the KEM simulatable and key indistinguishable, and the DDH assumption hold.

Then, the protocol in Fig. 4 securely realizes ideal functionality in Fig. 2 for the join defined in Def. 3 for semi-honest corruption of one of the two parties C , D and any amount of parties P_1 to P_T . In case of a corruption of D , the leakage graph of Def. 9 is leaked.

3.3 Rerandomizable Encrypted OPRF (EO)

OPRFs allow a client to obliviously evaluate a function PRF on their private input x with the server's secret key sk (i.e., $\text{PRF}_{sk}(x)$) [11]. We introduce a new rerandomizable encrypted OPRF (EO) primitive with more powerful functionality that allows: (a) multiple input providers to encrypt their inputs, (b) an output receiver to shuffle and rerandomize the ciphertexts, (c) a server to obliviously evaluate a PRF on encrypted as well as plaintext identifiers, and (d) the output receiver to decrypt the encrypted PRF evaluations. Our EO primitive consists of a collection of seven algorithms:

DEFINITION 11 (EO). A rerandomizable encrypted OPRF (EO) parameterized with security parameter κ is a collection of algorithms (KG, EKG, Eval, Enc, Rnd, OEval, Dec) with the following syntax.

- $\text{KG}(1^K)$: On input 1^K output a public key, secret key pair (pk, sk) .

- $\text{EKG}(1^K)$: On input 1^K output a public function key, evaluation key pair (pf, ek) .
- $\text{Eval}(ek, x)$: On input (ek, x) , output a PRF output y .
- $\text{Enc}(pk, pf, x)$: On input (pk, pf, x) , output a ciphertext ct .
- $\text{Rnd}(pk, pf, ct)$: On input (pk, pf, ct) , output a ciphertext ct' .
- $\text{OEval}(ek, ct)$: On input (ek, ct) , output evaluated ciphertext ect .
- $\text{Dec}(sk, ect)$: On input (sk, ect) , output y .

For correctness, we ask that for any $x \in \{0, 1\}^*$, $\Pr[\text{Dec}(sk, \text{OEval}(ek, \text{Rnd}(pk, pf, \text{Enc}(pk, pf, x)))) = \text{Eval}(ek, x)] \geq 1 - \text{negl}$, where $(pk, sk) \leftarrow \text{KG}(1^K)$ and $(pf, ek) \leftarrow \text{EKG}(1^K)$.

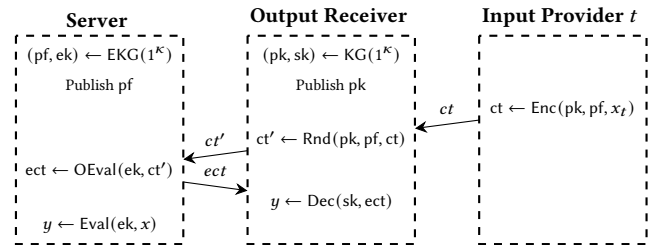


Figure 5: Rerandomizable Encrypted OPRF (EO).

We use EO as shown in Fig. 5. Each *Input Provider* t invokes EO.Enc to encrypt identifier x_t . Afterward, an *Output Receiver* rerandomizes and shuffles the ciphertexts using EO.Rnd . We remark that for security, we require that neither possession of EO.sk nor EO.ek is sufficient to distinguish encryptions of two different messages. After the shuffle, the *Server* uses EO.OEval to homomorphically evaluate a PRF on the encrypted identifier. The *Server* also evaluates the PRF on plaintext identifiers without knowledge of EO.sk by using EO.Eval and knowledge of EO.ek . Finally, the *Output Receiver* uses EO.Dec to decrypt the PRF evaluation. Observe that both the *Server* and the *Output Receiver* end up with the same PRF evaluation y for the same input x (or x_t). In order to have a PRF, we require that the EO.Eval outputs are pseudorandom given EO.pk , EO.sk , and EO.pf . In Appendix B, we define several security notions and show how to construct this primitive from DDH.

3.4 DPMC with Secure Shuffling ($D_s\text{PMC}$)

In DPMC, party D performs the left join on the hashed and exponentiated data between C and multiple delegators denoted as P_t . This process enables D to learn the full bipartite graph of correlations of matches up to an isomorphism due to shared identifiers. We address this issue with an enhanced version called $D_s\text{PMC}$ that utilizes our novel EO scheme and employs two delegates: party D and a new *shuffler* party S . $D_s\text{PMC}$ relies on EO to perform a secure three-party shuffling protocol between C , D , and S that combines and shuffles the data from all delegator parties P_t . In the process of secure shuffling, the data from the delegators are reordered in a way that no single party knows the applied permutation. Additionally, the delegators' data undergo two forms of rerandomization. First, the encrypted identifiers are refreshed with new ciphertexts using the EO.Rnd algorithm, generating new ciphertexts that correspond to the same plaintexts. Second, the secret shares of the associated data are reshared, creating new secret shares of the same plaintext

Setup: All parties P_t have access to the public key pk_D of party D , party D has secret key sk_D . $M := \sum_{t=1}^T m_t$.

<p>① Key-Generation (Party C)</p> <p>1: $(\text{KEM.pk}, \text{KEM.sk}) \leftarrow \text{KEM.KG}(1^\kappa)$</p> <p>2: $(\text{EO.pf}, \text{EO.ek}) \leftarrow \text{EO.EKG}(1^\kappa)$</p> <p>Send to P_t, S, D: $\text{KEM.pk}, \text{EO.pf}$</p>	<p>⑦ Mask Shares (Party S)</p> <p>Input: $\{\widetilde{\text{sh}}_{C,i}\}_{i \in [M]}$</p> <p>Messages: KEM.pk</p> <p>1: For i in $[M]$: ▷ For all delegators rows.</p> <p>2: $(\text{KEM.cp}_i, \text{KEM.k}_i) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ ▷ New encaps.</p> <p>3: $\widetilde{\text{sh}}_{D,i} := \text{sh}_{C,i} \oplus \text{KEM.k}_i$ ▷ Updated shares for party D.</p> <p>Send to D: $\{\text{KEM.cp}_i, \widetilde{\text{sh}}_{D,i}\}_{i \in [M]}$</p>
<p>② Key-Generation (Party D)</p> <p>1: $(\text{EO.pk}, \text{EO.sk}) \leftarrow \text{EO.KG}(1^\kappa)$</p> <p>Send to P_t: EO.pk</p>	<p>⑧ Prepare Match Keys (Party C)</p> <p>Input: $\text{KV}_C = \{(c_{i,j})_{j \in [m_C,i]}\}_{i \in [m_C]}$, EO.ek and $\{\widetilde{\text{EO.ct}}_{i,j}\}_{i \in [M], j \in [m_i]}$.</p> <p>1: For i in $[M]$, j in $[m_i]$: ▷ For all delegators rows and columns.</p> <p>2: $\text{EO.ect}_{i,j} := \text{EO.OEval}(\text{EO.ek}, \widetilde{\text{EO.ct}}_{i,j})$</p> <p>3: For i in $[m_C]$, j in $[m_C,i]$: ▷ For all rows and columns in KV_C.</p> <p>4: $\text{hc}_{i,j} := \text{EO.Eval}(\text{EO.ek}, c_{i,j})$</p> <p>5: Use $c_i := (c_{i,j})_{j \in [m_C,i]}$ to order $(\text{hc}_{i,j})_{j \in [m_C,i]}$.</p> <p>Send to D: $(\text{hc}_{i,j})_{i \in [m_C], j \in [m_C,i]}$, $\{\text{EO.ect}_{i,j}\}_{i \in [M], j \in [m_i]}$</p>
<p>③ Identity Match (Party P_t)</p> <p>Input: $\text{KV}_t = \{(p_{t,i,j})_{j \in [m_t,i]} v_{t,i}\}_{i \in [m_t]}$ for data set size m_t.</p> <p>Messages: $\text{EO.pk}, \text{EO.pf}$</p> <p>1: $\text{seed}_t \xleftarrow{R} \{0,1\}^\kappa$ ▷ Random seed t.</p> <p>2: $\text{cta}_t := \text{PKE.Enc}(\text{pk}_D, \text{seed}_t)$</p> <p>3: $(\text{sh}_{D,t,1}, \dots, \text{sh}_{D,t,m_t}) \xleftarrow{R} \text{PRG}(\text{seed}_t)$ ▷ Share of $v_{t,i}$ for party D.</p> <p>4: For i in $[m_t]$: ▷ For each row in KV_t.</p> <p>5: For j in $[m_t,i]$: ▷ For each column.</p> <p>6: $\text{EO.ct}_{t,i,j} \leftarrow \text{EO.Enc}(\text{EO.pk}, \text{EO.pf}, p_{t,i,j})$ ▷ Encrypt data using EO.</p> <p>7: $\text{sh}_{C,t,i} := v_{t,i} \oplus \text{sh}_{D,t,i}$ ▷ Share of $v_{t,i}$ for party C.</p> <p>Send to C: $\text{cta}_t, \{(\text{EO.ct}_{t,i,j})_{j \in [m_t,i]}, \text{sh}_{C,t,i}\}_{i \in [m_t]}$</p>	<p>⑨ Identity Match and Recover Shares (Party D)</p> <p>Input: EO.sk and $\{\widetilde{\text{sh}}_{D,i}\}_{i \in [M]}$</p> <p>Messages: KEM.pk, $\{\text{KEM.cp}_i, \widetilde{\text{sh}}_{D,i}\}_{i \in [M]}$, $\{\text{hc}_{i,j}\}_{i \in [m_C], j \in [m_C,i]}$, $\{\text{EO.ect}_{i,j}\}_{i \in [M], j \in [m_i]}$</p> <p>1: For i in $[M]$, j in $[m_i]$: ▷ For all delegators rows and columns.</p> <p>2: $\text{h}_{i,j} := \text{EO.Dec}(\text{EO.sk}, \text{EO.ect}_{i,j})$</p> <p>3: $\mathcal{J} := (\text{hc}_{i,j})_{i \in [m_C], j \in [m_C,i]} \triangleright (\text{hc}_{i,j})_{i \in [M], j \in [m_i]}$ ▷ Details in Alg. 2.</p> <p>4: For each row i in \mathcal{J}, repeat for $t \in [T]$: ▷ For each row in the join.</p> <p>5: If record matched:</p> <p>6: $\widetilde{\text{KEM.cp}}_i := \text{KEM.cp}_i$ ▷ Use encaps. from P_t.</p> <p>7: $\widehat{\text{sh}}_{D,i,t} := \widetilde{\text{sh}}_{D,i'} \oplus \widetilde{\text{sh}}_{D,i'}$ ▷ Final shares for party D.</p> <p>8: Else: ▷ no match found</p> <p>9: $(\widetilde{\text{KEM.cp}}_i, \text{KEM.k}_i) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ ▷ New encaps.</p> <p>10: $\widehat{\text{sh}}_{D,i,t} := \text{KEM.k}_i$ ▷ Use share of 0.</p> <p>11: Pick m_C random permutations $\{\pi_i\}_{i \in [m_C]}$.</p> <p>12: $\mathcal{J}_D := (\pi_i(\{\widehat{\text{sh}}_{D,i,t}\}_{t \in [T]}))_{i \in [m_C]}$ ▷ D's permuted XOR shares.</p> <p>Send to C: $\{\pi_i(\{\text{KEM.cp}_{i,t}\}_{t \in [T]})\}_{i \in [m_C]}$</p>
<p>④ Forward Shares to D (Party C)</p> <p>Messages: $\{(\text{EO.ct}_{t,i,j})_{j \in [m_t,i]}, \text{sh}_{C,t,i}\}_{i \in [m_t]}, \text{cta}_t, \}_{t \in [T]}$</p> <p>Send to D: $\{\text{cta}_t\}_{t \in [T]}$</p>	<p>⑩ Recover Shares (Party C)</p> <p>Input: KEM.sk</p> <p>Messages: $\{\widetilde{\text{KEM.cp}}_i\}_{i \in [m_C], t \in [T]}$</p> <p>1: For i in $[m_C]$, t in $[T]$</p> <p>2: $\widehat{\text{sh}}_{C,i,t} := \text{KEM.Dec}(\text{KEM.sk}, \widetilde{\text{KEM.cp}}_i)$ ▷ $(\text{KEM.k}_{i,t})$ Shares for C.</p> <p>3: $\mathcal{J}_C := (\widehat{\text{sh}}_{C,i,t})_{i \in [m_C], t \in [T]}$ ▷ Aligned with $(c_i)_{i \in [m_C]}$</p>
<p>⑤ Reconstruct Shares (Party D)</p> <p>Messages: $\{\text{cta}_t\}$ for all T parties P</p> <p>1: For t in $[T]$: ▷ For each delegator P_t</p> <p>2: $\text{seed}_t := \text{PKE.Dec}(\text{sk}_D, \text{cta}_t)$ ▷ Get seed seed_t.</p> <p>3: $\text{sh}_{D,t,1}, \dots, \text{sh}_{D,t,m_t} := \text{PRG}(\text{seed}_t)$ ▷ Share of $v_{t,i}$ for party D.</p>	<p>⑥ Shuffling – Appendix C (Parties C, S, D)</p> <p>Note: The EO.ct ciphertexts as well as the sh_C and sh_D secret shares are: a) reordered such that no party knows the permutation, and b) rerandomized to $\widetilde{\text{EO.ct}}$, $\widetilde{\text{sh}}_C$, and $\widetilde{\text{sh}}_D$. These rerandomizations correspond to fresh encryptions and fresh secret shares of the same underlying data.</p> <p>C Input: $\{(\text{EO.ct}_{t,i,j})_{j \in [m_t,i]}, \text{sh}_{C,t,i}\}_{i \in [m_t], t \in [T]}$</p> <p>S Input: –</p> <p>D Input: $\{\text{sh}_{D,t,i}\}_{i \in [m_t], t \in [T]}$</p> <p>C receives output: $\{(\widetilde{\text{EO.ct}}_{i,j})_{i \in [M], j \in [m_i]}\}$ ▷ Rerandomized ciphertexts.</p> <p>S receives output: $\{\widetilde{\text{sh}}_{C,i}\}_{i \in [M]}$ ▷ Rerandomized shares for party C.</p> <p>D receives output: $\{\widetilde{\text{sh}}_{D,i}\}_{i \in [M]}$ ▷ Rerandomized shares for party D.</p>

Figure 6: Multi-key D_S PMC. This protocol uses two delegate parties (S and D) and is based on secure shuffling and EO.

values. Notably, these rerandomization steps do not reveal the underlying data and are meant to break any link between the data that the delegators provide and the data that are used for the join and the secret sharing. This way, the leakage to party D is only between C 's data and the combined data of parties P_1 to P_T , contrary to the pairwise leakages of DPMC. Since C combines the inputs of all delegators, party D (who performs the join) only sees two encrypted datasets (i.e., encrypted KV_C and encrypted KV_D).

Due to our shuffling and rerandomization steps, in a potential corruption of the delegators and one of C, D, S , the corrupted parties cannot infer any information as the data have been permuted and rerandomized. Our shuffling scheme is secure in the honest-majority setting, which is the case with multiple applications from

both academia [1, 15, 39] and industry. For instance, Mozilla recently deployed a service that relies on the Prio protocol to collect telemetry data about Firefox [30], while Crypten [33] and TF Encrypted [17] build privacy-preserving machine learning frameworks for PyTorch and TensorFlow, respectively. We delve into the details of the security of D_S PMC in Appendix D.2.

D_S PMC follows a similar approach as DPMC, with the difference that it leverages our EO primitive. We formally present our D_S PMC protocol in Fig. 6; intuitively, it works as follows. The delegator parties use EO to encrypt their identifiers and generate XOR shares as in the DPMC protocol. Then, the delegators encrypt the shares for party D using pk_D and send them to Party C along with C 's shares and all of the ciphertexts. C then forwards D 's encrypted

shares to Party D who decrypts them. Then parties C , D , and S run a secure shuffling protocol in which C receives rerandomized EO ciphertexts ($\widehat{\text{EO.ct}}$), S obtains the rerandomized shares for C ($\widehat{\text{sh}}_C$), and D receives its randomized shares ($\widehat{\text{sh}}_D$).

Next, S generates new key encapsulations and uses C 's shares to generate new shares for D . It sends the updated $\widehat{\text{sh}}_{D,i}$ to party D that allows D to adjust their shares to be consistent with the new shares of C . S also sends the encapsulations ($\widehat{\text{KEM.cp}}$) to D .

Party C proceeds by homomorphically evaluating the PRF on the EO ciphertexts and the PRF on its own identifiers $c_{i,j}$ and sends the outcomes to D . Recall from Fig. 5 that all the output receiver (D in this case) needs to do now is to decrypt the evaluated EO ciphertexts and compute the matches. D computes the left join with associated data, where for each matched record it ends up with T shares of either the delegators' associated data or zero. Note that D does not know which decrypted PRF identifier belongs to which delegator. The matching logic is described in more detail in Def. 4 in the Appendix A. Similar to DPMC, D it replaces the encapsulation with a fresh encapsulation when no match is found and keeps $\widehat{\text{KEM.k}}$ as its share $\widehat{\text{sh}}_D$. It then forwards the encapsulations $\widehat{\text{KEM.cp}}$ to C . Party C finalizes the protocol by recovering the encapsulated keys and using them as its shares $\widehat{\text{sh}}_C$. Observe that for the unmatched records, C and D end up with secret shares of zero as $\mathcal{J}_C \oplus \mathcal{J}_D = \widehat{\text{sh}}_{C,i,\hat{i}} \oplus \widehat{\text{sh}}_{D,i,\hat{i}} = \widehat{\text{KEM.k}}_{i,\hat{i}} \oplus \widehat{\text{KEM.k}}_{i,\hat{i}} = 0$, while for the matched records we get secret shares of the delegators' associated data as $\mathcal{J}_C \oplus \mathcal{J}_D = \widehat{\text{sh}}_{C,i,\hat{i}} \oplus \widehat{\text{sh}}_{D,i,\hat{i}} = \widehat{\text{KEM.k}}_i \oplus \widehat{\text{sh}}_{D,i} \oplus \widehat{\text{sh}}_{D,i} = \widehat{\text{KEM.k}}_i \oplus (\widehat{\text{sh}}_{D,i} \oplus \widehat{\text{sh}}_{C,i}) \oplus \widehat{\text{KEM.k}}_i = \widehat{\text{sh}}_{D,i} \oplus \widehat{\text{sh}}_{C,i} = \widehat{\text{sh}}_{D,i} \oplus \widehat{\text{sh}}_{C,i} = v_i$.

We describe our protocol's leakage in Def. 12. $D_s\text{PMC}$ limits the leakage of DPMC from pairwise intersection sizes between each party to one intersection size between party C and the union of all delegators. For instance, for parties C , P_1 , and P_2 , party D will learn $|\text{KV}_C \cap \text{KV}_P|$, where $\text{KV}_P := \{\text{KV}_1 \cup \text{KV}_2\}$. Notably, these intersection sizes also contain the number of times that keys are matched (i.e., 1 to T). In case multiple keys are used, D will additionally learn a graph of matches as defined by $L_{x,y}$ in Def. 12. We provide the security of $D_s\text{PMC}$ in Theorem 13 and prove it in Appendix D.2. Note that we do not need ciphertext indistinguishability for the secret key owner (Lemma 30) since D does not handle any EO ciphertexts, only evaluated ciphertext. This might change when a different shuffle protocol is used.

DEFINITION 12 ($D_s\text{PMC}$ LEAKAGE). *Given KV_C and $\text{KV}_1, \dots, \text{KV}_T$, the leakage $L_{x,y}$ of the ideal functionality in Fig. 2 for the $D_s\text{PMC}$ protocol in Fig. 6 is defined as follows. Merge $\text{KV}_1, \dots, \text{KV}_T$ to $\text{KV}_P := \bigcup_{t \in [T]} \text{KV}_t$. Define $\text{KV}_{u,C}$ by replacing $c_{i,j} \in \text{KV}_C$ with $u_{i,j} \xleftarrow{\mathcal{R}} \{0, 1\}^K$. Define $\text{KV}_{u,P}$ by replacing $p_{i,j} \in \text{KV}_P$ with $u_{i',j'}$ if there exists an i', j' pair with $p_{i,j} = c_{i',j'}$ or an already replaced $p_{i',j'}$ with $p_{i,j} = p_{i',j'}$, otherwise replace it with $u'_{i,j} \xleftarrow{\mathcal{R}} \{0, 1\}^K$. $L_{x,y} := \{(C, \text{KV}_{u,C}), (D, \text{KV}_{u,P})\}$.*

THEOREM 13. *Let PKE be an IND-CPA secure and correct PKE scheme, KEM a correct and key-indistinguishable key encapsulation mechanism, PRG as secure pseudorandom generator, and EO be a correct and satisfy statistical rerandomized ciphertext indistinguishability, the (semi-honest) ciphertext indistinguishability for the evaluation key and secret key owner and ciphertext well-formedness.*

Then, the protocol in Fig. 6 securely realizes ideal functionality in Fig. 2 for the join defined in Def. 4 for semi-honest corruption of one of the three parties C , D , S , and any amount of parties P_1 to P_T . In case of a corruption of D , the leakage graph of Def. 12 is leaked.

4 MATCHING STRATEGY

Recall from Fig. 1 that the view of each party for a specific record may be different and a record may have multiple identifiers (e.g., email address, phone). When combining datasets from multiple delegators, the uniqueness of the identifiers cannot be guaranteed as the same record might appear in more than one dataset. Thus, potential matches for each row can occur based on different identifiers across different delegators. For instance, a match on the j th identifier of record c_i may occur for keys in different positions between different parties (e.g., with $p_{t,i',j'}$ with $i \neq i'$ and $j \neq j'$). Parties C and D in our protocols compute the left join as described in Def. 3 and acquire \mathcal{J}_C and \mathcal{J}_D , respectively. To capture all the aforementioned matches, for T delegators, \mathcal{J}_C and \mathcal{J}_D have T permuted columns of secret shares which either correspond to shares of the associated metadata of one of the input parties (if a match was found) or to shares of NULL (in case no match was found).

As the number of delegators P_1 to P_T grows, it is natural for our resulting \mathcal{J}_C and \mathcal{J}_D tables to contain multiple secret shares of NULL. This becomes more evident if each individual dataset KV_t is relatively small compared to KV_C ; even if all the records of KV_t match with records in KV_C , there would still be multiple unmatched records in KV_C which will get secret shares of NULL. To optimize both our matching and our downstream computation, we now delve into a matching strategy to generate one-to-many connections that do not depend on T and minimize the number of NULL secret shares.

First, C and D agree on a maximum number of connections K to capture. D performs a ranked left join by starting from the identifier with the highest priority in KV_C and checking whether it appears in each KV_t before moving to the next record in KV_C . After searching by the first key of each record in KV_C , D continues with the next identifier, and so on. If a record from P_t is matched, we mark that record as done and continue to the next record in order to avoid counting the same associated values more than once. For each record c_i , if K or more matches are found, D creates secret shares of the associated data of the first K records, otherwise (if less than K matches are found), D pads the remaining columns (up to K) with secret shares of NULL. We note that this is an implementation-specific detail that can be trivially extended to different matching strategies. Each of the resulting tables \mathcal{J}_C and \mathcal{J}_D has K columns and captures a one-to- K matches for each record in the left join.

5 REAL-WORLD APPLICATIONS

Recall that our ideal functionality $\mathcal{F}_{\text{DPMC}}$ (and $\mathcal{F}_{\text{DPMC}}$) consists of $\mathcal{F}_{\text{JOIN}}$ and \mathcal{F}_{CMP} . Our delegated protocols realize $\mathcal{F}_{\text{JOIN}}$ and output secret shares to parties C and D for the left join of parties C and P_1, \dots, P_T . Next, \mathcal{F}_{CMP} can be realized by running any general-purpose MPC between C and D . We foresee multiple real-world applications for \mathcal{F}_{CMP} that may leverage our architecture merging multiple private datasets across distrusting parties with a centralized entity (party C) to securely compute analytics. For instance, DPMC enables calculating the risk of a health condition by merging

information held by a larger healthcare provider with data stored on millions of individual smart devices. In another example, an ad publisher holding user-provided information can measure advertising efficacy and offer personalization by merging with data held by multiple advertisers while still preserving user privacy. In this section, we focus on the latter and outline how DPMC enables privacy-preserving ad measurement and delivery of personalized advertising leveraging privacy-preserving machine learning. The former provides advertisers useful insights about how their ad campaigns are performing, while the latter enables delivering personalized ads while preserving user privacy.

5.1 Privacy-Preserving Ad Attribution

Inputs. We assume the following input data held by an ad publisher, denoted by C and T advertisers, denoted by P_1, \dots, P_T .

- *Party C* is a company that holds a dataset of ad actions (i.e., clicks) performed by individuals on product-related advertisements. These ads were shown to users after they expressed an intent via an online search engine. Users may be shown ads related to multiple products owned by hundreds of advertisers.
- Advertisers P_1 to P_T , hold conversion information for their customers, such as purchase amount and time of the purchase.
- All parties (C, P_1, \dots, P_T) also hold annotated sets of common identifiers (e.g., email addresses and phone numbers).

$\mathcal{F}_{\text{JOIN}}$ phase. Executing the DPMC protocol for $\mathcal{F}_{\text{JOIN}}$ with the above input data from C and multiple P parties, the following output is available at the ad publisher C and the delegate servers.

- *Party C* holds a mapping of secret shares of conversion data to a dataset of ad actions. This mapping does not reveal any new information to C apart from random-looking secret shares. In the case of no matches, party C receives secret shares of zero.
- *Party D* receives a set of secret shares of the conversion data or a dummy value (e.g., zero) that is also aligned to party’s C records (i.e., left join). D gains insights into pairwise intersection sizes (in DPMC) or the intersection size of C with the union of all advertisers’ sets (in $D_s\text{PMC}$). For example, when users have unique phone numbers and email addresses, in $D_s\text{PMC}$ D learns the intersection sizes of records where at least phone number, email address, or both matched between the company and the union of all advertisers’ data. In a real-world scenario where D is a privacy-conscious non-profit organization, this level of leakage has fairly low privacy implications. If the uniqueness of identifiers cannot be assumed, the sizes of groups with the same identifiers are leaked.
- Parties P_1 to P_T , receive nothing.

\mathcal{F}_{CMP} phase. Parties C and D now hold secret shares of conversion metadata such as conversion time and values. C can then further input metadata of ad actions, such as click timestamp, as secret shares using the link to the original records that were established by the DPMC protocol. Now, parties C and D engage in multi-party computation to compute the attribution function that flags when a conversion (product was bought) occurred within a pre-specified time window from the ad action. Note that the MPC computation is embarrassingly parallel given the row-wise output structure of DPMC. The output of the privacy-preserving ad attribution remains

at the ad action level, hence remains secret shared between parties C and D and is used as an input into further downstream computations such as private measurement or personalization, described next.

5.2 Privacy-Preserving Analytics

Measuring the efficacy of advertising first requires computing aggregated conversion outcomes such as the total number of attributed conversions per campaign. Note that DPMC maintains the left join of the ad actions without revealing any user-level information to party C at any stage. Party C may attach campaign-level identifiers with limited entropy ensuring sufficient K -anonymity guarantees. At this point, parties C and D engage in another round of MPC (i.e., a new \mathcal{F}_{CMP} phase) to compute aggregated conversion outcomes per campaign. Finally, differentially private noise can be added to the aggregated outcomes within MPC before revealing the results to party C , so that C only learns noisy aggregates for each ad campaign.

5.3 Privacy-Preserving Personalization

Privacy-preserving personalization typically entails training a model to be able to estimate the relevance of potential ads for users. Note that privacy-preserving ad attribution during the data pre-processing phase generates secret shares of ad attribution outcomes for both parties C and D . Leveraging the mapping produced by DPMC from secret shares to original ad actions, party C may attach any private features to the private attribution outcomes without revealing any individually identifiable information. At this stage, parties C and D can run a new \mathcal{F}_{CMP} in multi-party computation for model training with privately input features (from party C) and secret shared labels (from both parties C and D). For example, CrypTen [33], a multi-party computation framework for machine learning, may be leveraged between the parties C and D downstream to the DPMC protocol. Similarly to the aforementioned analytics example, privacy-preserving personalization would also include differential privacy guarantees and we point avid readers to one such implementation [55].

6 EVALUATIONS

Implementation & Setup. We implemented our protocols in Rust (1.62) and used the Dalek library for Elliptic Curve Cryptography with Ristretto for Curve25519 [18, 28].² This enables the use of a fast curve while avoiding high-cofactor vulnerabilities. For symmetric encryption, we use the Fernet library with AES-128 in CBC mode, for public key encryption we use ElGamal with elliptic curves, and for the key encapsulation mechanism, we use ElGamal KEM.

We created artificial datasets where each record has one 128-bit identifier and two 64-bit associated values. The performance measurements were carried out on AWS m5.12xlarge EC2 instances (Intel Xeon at 3.1GHz, 48 vCPU, 192GB RAM). To simulate C , D , and multiple P parties we leverage three separate EC2 instances in the same region, where C and D are hosted by two separate instances, and the third instance hosts all parties P_1 to P_T . For our WAN experiments, we used three m5.12xlarge EC2 instances in N. Virginia, Ohio, and N. California. All parties communicate via RPC over TLS v1.3 using Protocol Buffers.

²Our protocols are open-source at <https://github.com/facebookresearch/Private-ID>.

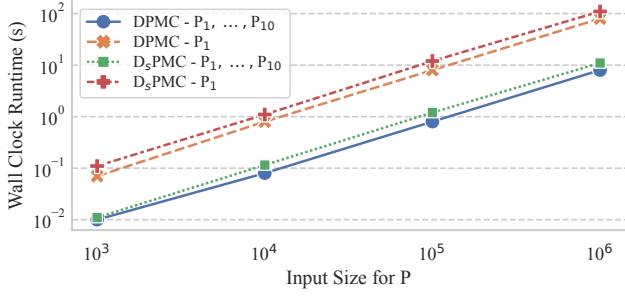


Figure 7: Measured execution time of delegator for DPMC and D_s PMC with $m_C = 10^6$ and increasing KV_P with intersection sizes of 50% of KV_P .

Varying number of delegators. In Fig. 7, we fixed the size of KV_C to 1 million and varied both the number of delegators and their dataset sizes. In the orange and red trends, we used a single delegator for DPMC and D_s PMC with dataset sizes indicated by the x-axis. In the blue and green trends, we split the dataset into ten delegators, where each party has 1/10 of the input size shown on the x-axis. Although the combined size of the dataset of the ten parties is the same, the local computation for each delegator is significantly less. In this case, the performance time for each delegator is about ten times faster than having a single P_1 party with a bigger dataset.

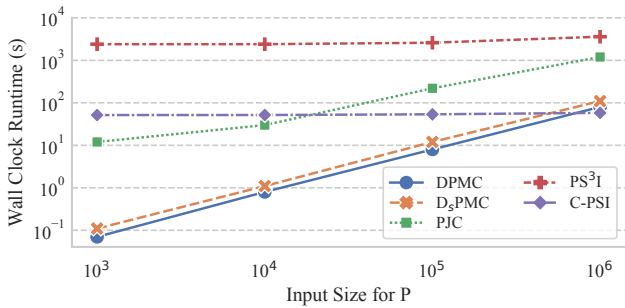


Figure 8: Measured time of P for DPMC, D_s PMC, PJC, PS^3I , and Circuit-PSI with $m_C = 10^6$ and intersection sizes of 50% of m_t . All protocols are evaluated with a single delegator.

Protocol time for delegator. Next, we fixed the input of party C to 1 million with a single identifier per record and varied the size of the dataset of the delegator. In Fig. 8 we show the execution times for party P for DPMC, D_s PMC, PJC, PS^3I , and Circuit-PSI. We use the PS^3I [10] and PJC [31] implementations from [9], which both use Paillier with a 2048-bit public key. Similarly to our protocols, both these protocols assume that party P has associated metadata: PS^3I generates additive secret shares, whereas PJC aggregates the associated values of the items in the intersection. Additionally, for fair comparisons, we implemented over-the-network communication between the sender and the receiver on the Circuit-PSI implementation of [49] (the implementation of [12] crashes with different dataset sizes). The blue and orange trends in Fig. 8 show the execution time of a single delegator running the DPMC and the

D_s PMC protocols, respectively, which are approximately the same. We observe that the execution time for both is approximately $10\times$ faster than party P in PJC and multiple orders of magnitude faster than PS^3I . The runtime in PJC scales linearly with P 's dataset size, however, this is not the case with PS^3I as P 's execution time is also affected by C 's dataset. The runtime for Circuit-PSI is linear in the input of both parties, so it incurs high overheads when $m_C \gg m_P$. Both this and the previous experiments (Figs. 7 and 8) demonstrate the benefits of our delegated protocols for the delegator parties compared to the two-party protocols.

Varying intersection size. In our next experiment, we fixed the input size at 1 million records for both parties and the number of identifiers at 2 per record and varied the intersection size (1%, 25%, 50%, and 100%). We observed negligible performance variations (i.e., less than a second) for the different intersection sizes since our protocol always outputs the left join and depends on KV_C size.

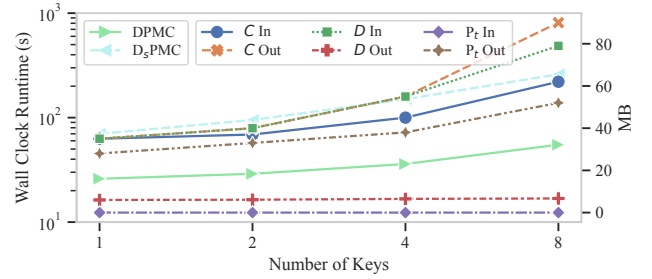


Figure 9: Wall clock time for DPMC and D_s PMC. DPMC network traffic for each party C , D , and P_t with an increasing number of keys per row for $m_C = m_t = 10^5$ and an intersection size of 50% of m_C .

Varying number of identifiers. We now show how the number of keys affects the performance of our protocols. We fixed the input size to 10^5 records for both parties and the intersection size to 50%. Fig. 9 shows the total time for DPMC (light green trend) and D_s PMC (light blue trend) as well as the input and output traffic for each party for DPMC. Notably, the communication of D_s PMC is similar for an increasing number of identifiers as the matching strategy is very similar for the two protocols.

Table 2: Communication cost & number of exponentiation. T is the number of delegators; m_C and m_t are the set sizes of C and each delegator P_t , respectively, and $M := \sum_{t=1}^T m_t$. $\mathcal{I} := |KV_C \cap KV_P|$, where $KV_P := \{KV_1 \cup KV_2 \cup \dots \cup KV_T\}$.

	Party	C	D	S	P_t
DPMC	Communication	$\mathcal{O}(m_C + M)$	$\mathcal{O}(m_C + M)$	-	$\mathcal{O}(m_t)$
	Num. of Exp.	$2m_C + M$	$M + m_C - \mathcal{I}$	-	$2m_t + 1$
D_s PMC	Communication	$\mathcal{O}(m_C + M)$	$\mathcal{O}(m_C + M)$	$\mathcal{O}(M)$	$\mathcal{O}(m_t)$
	Num. of Exp.	$2m_C + 3M$	$M + m_C - \mathcal{I}$	$4M$	$2m_t$

Communication. In Table 2 we present the asymptotic costs (communication and number of exponentiations) of each protocol for each party. C and D incur similar communication overhead which

scales with the size of KV_C and the delegators' datasets. The communication cost for each delegator P_t is linear to their dataset. In D_S PMC, party S incurs a linear communication to the size of all the delegators' datasets. Finally, we observe that the number of exponentiations of DPMC and D_S PMC are similar.

Table 3 shows each party's incoming and outgoing traffic in MBs. We observe a linear increase in the communication for each party as we increase the input sizes. Interestingly, we see that although D_S PMC performs more rounds than DPMC, the communication for each party is lower than DPMC. This happens because each P_t encrypts their XOR shares in order to prevent C from accessing them during the fourth step of the protocol. Finally, our protocols have similar communication as Circuit-PSI, which showed a linear increase with the dataset sizes. For reference, the outgoing communication for datasets of 10^6 elements was 424 MBs (344 from the sender and 80 from the receiver).

Table 3: For each party C, P, D, S we show In/Out in MB with $m_C = m_P$ and intersection size $I = 50\%$ of m_C .

	Size	C [In/Out]	D [In/Out]	S [In/Out]	P_t [In/Out]
DPMC	10^3	0.3/0.3	0.3/0.1	-	0.1/0.3
	10^4	3.7/3.7	3.4/2.8	-	0.1/2.8
	10^5	33/33	33/4.8	-	0.1/28
	10^6	312/312	320/44	-	0.1/279
D_S PMC	10^3	0.2/0.2	0.1/0.1	0.1/0.1	0.1/0.1
	10^4	2.3/2.5	1.5/0.4	1/1	0.1/0.8
	10^5	22/24	14/4.3	9.5/9.5	0.1/8.5
	10^6	220/241	145/42	94/94	0.1/84.7

Two-party related works. We also compare our protocols with two-party works and vary the input size of each party from 10^3 to 10^6 while fixing the intersection size to 50%. As a baseline, we compare with multi-key Private-ID [8] which only focuses on private matching and does not consider associated data. To be in a similar setting, we run our protocols with a single party P . Fig. 10 shows how our delegated protocols significantly outperform both PS^3I and PJC by more than a factor of 10x. On the other hand, our delegated protocols are only $\approx 1.8x$ slower than Private-ID although the latter does not include associated values and is only between two parties. This means that our protocols process approximately twice the amount of data that Private-ID processes since for each row of KV_P , DPMC and D_S PMC also create secret shares of the associated data. Circuit-PSI is 3-4x times faster than our protocols but, as the other related works, only considers two parties who are both assumed to be online throughout the entire protocol execution and do not take into account any delegation methods.

In Fig. 11, we repeated the same experiments over WAN and observed a similar scaling for all the protocols. Interestingly, the margin between Circuit-PSI and our protocols became smaller as Circuit-PSI requires significantly more communication. Finally, note that these experiments do not offer a balanced assessment of our protocols as the benefits of the delegated setting are shown in Figs. 7 and 8. In the former, we observe that each delegator performs work proportional to their dataset size, while in the latter, the delegators have smaller datasets than C and they go offline after they outsource their datasets.

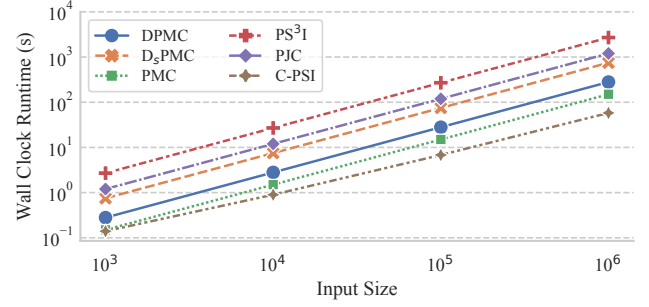


Figure 10: Comparisons of DPMC and D_S PMC with two-party protocols: PJC, Private-ID (PMC), PJC, PS^3I , and Circuit-PSI with an increasing number of dataset sizes ($m_C = m_P$) and an intersection of 50% m_C . We use PMC as a baseline as it only performs matching and does not consider associated values.

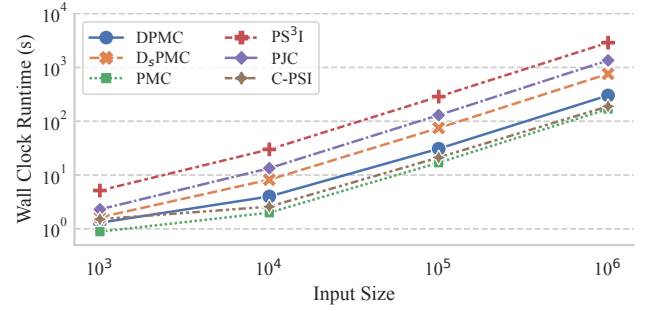


Figure 11: Comparisons as in Fig. 10 over WAN.

7 CONCLUDING REMARKS

We presented two delegated protocols that establish relations between datasets that are held by multiple distrusting parties and enable them to run any arbitrary secure computation. Our protocols allow the input parties to submit their records along with associated values and generate secret shares of the associated values for the matched records and secret shares of NULL otherwise. Notably, they facilitate the delegation of both the matching process and downstream secure computation to delegate parties. In contrast with prior works that only support two parties, our work is designed to scale to multiple input parties.

In addition, our delegated protocols enable one of the input parties to provide more data after the matching has been established which can be used for the downstream computation without requiring rerunning the private matching process. We further introduced a rerandomizable encrypted OPRF (EO) primitive that extends beyond the classic two-party OPRF setting and allows multiple input providers to interact with an output receiver and a server and perform oblivious PRF evaluations. While prior works mostly focused on intersection and union, we focused on left-join matching and we demonstrated its benefits in privacy-preserving online advertising by performing private ad attribution measurement, privacy-preserving analytics, and personalization. Finally, our implementation demonstrates the efficiency of our constructions by outperforming related works.

ACKNOWLEDGMENTS

The authors would like to thank Anderson Nascimento, Erik Taubeneck, Gaven Watson, Sanjay Saravanan, Shripad Gade, Pratik Sarkar, and Charles Gouert for the fruitful discussions and the anonymous reviewers for their feedback. The third author was partially supported by NSF awards #2101052, #2200161, #2115075, and ARPA-H SP4701-23-C-0074.

REFERENCES

- [1] Toshinori Araki, Jun Furukawa, Yehuda Lindell, Ariel Nof, and Kazuma Ohara. 2016. High-Throughput Semi-Honest Secure Three-Party Computation with an Honest Majority. In *ACM CCS 2016: 23rd Conference on Computer and Communications Security*, Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi (Eds.). ACM Press, Vienna, Austria, 805–817. <https://doi.org/10.1145/2976749.2978331>
- [2] Amos Beimel. 2011. Secret-Sharing Schemes: A Survey. In *International Conference on Coding and Cryptology*. Springer, Berlin, Heidelberg, 11–46.
- [3] Mihir Bellare and Phillip Rogaway. 1993. Random Oracles are Practical: A Paradigm for Designing Efficient Protocols. In *ACM CCS 93: 1st Conference on Computer and Communications Security*, Dorothy E. Denning, Raymond Pyle, Ravi Ganesan, Ravi S. Sandhu, and Victoria Ashby (Eds.). ACM Press, Fairfax, Virginia, USA, 62–73. <https://doi.org/10.1145/168588.168596>
- [4] Abhishek Bhowmick, Dan Boneh, Steve Myers, Kunal Talwar, and Karl Tarbe. 2021. The Apple PSI system.
- [5] Erik-Oliver Blass and Florian Kerschbaum. 2023. Private Collaborative Data Cleaning via Non-Equi PSI. In *2023 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, San Francisco, CA, USA, 1419–1434. <https://doi.org/10.1109/SP46215.2023.10179396>
- [6] Dan Bogdanov, Sven Laur, and Jan Willemson. 2008. Sharemind: A Framework for Fast Privacy-Preserving Computations. In *ESORICS 2008: 13th European Symposium on Research in Computer Security (Lecture Notes in Computer Science, Vol. 5283)*, Sushil Jajodia and Javier López (Eds.). Springer, Heidelberg, Germany, Málaga, Spain, 192–206. https://doi.org/10.1007/978-3-540-88313-5_13
- [7] Dan Boneh, Elette Boyle, Henry Corrigan-Gibbs, Niv Gilboa, and Yuval Ishai. 2021. Lightweight Techniques for Private Heavy Hitters. In *2021 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, San Francisco, CA, USA, 762–776. <https://doi.org/10.1109/SP40001.2021.00048>
- [8] Prasad Buddharapu, Benjamin M Case, Logan Gore, Andrew Knox, Payman Mohassel, Shubho Sengupta, Erik Taubeneck, and Min Xue. 2021. Multi-key Private Matching for Compute. Cryptology ePrint Archive, Report 2021/770. <https://eprint.iacr.org/2021/770>.
- [9] Prasad Buddharapu, Andrew Knox, Payman Mohassel, Shubho Sengupta, Erik Taubeneck, and Vlad Vlaskin. 2020. Private-ID. <https://github.com/facebookresearch/Private-ID>.
- [10] Prasad Buddharapu, Andrew Knox, Payman Mohassel, Shubho Sengupta, Erik Taubeneck, and Vlad Vlaskin. 2020. Private Matching for Compute. Cryptology ePrint Archive, Report 2020/599. <https://eprint.iacr.org/2020/599>.
- [11] Silvia Casacuberta, Julia Hesse, and Anja Lehmann. 2022. SoK: Oblivious Pseudorandom Functions. Cryptology ePrint Archive, Report 2022/302. <https://eprint.iacr.org/2022/302>.
- [12] Nishanth Chandran, Divya Gupta, and Akash Shah. 2022. Circuit-PSI With Linear Complexity via Relaxed Batch OPPRF. *Proceedings on Privacy Enhancing Technologies* 2022, 1 (Jan. 2022), 353–372. <https://doi.org/10.2478/popets-2022-0018>
- [13] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. 2018. Labeled PSI from Fully Homomorphic Encryption with Malicious Security. In *ACM CCS 2018: 25th Conference on Computer and Communications Security*, David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang (Eds.). ACM Press, Toronto, ON, Canada, 1223–1237. <https://doi.org/10.1145/3243734.3243836>
- [14] Michele Ciampi and Claudio Orlandi. 2018. Combining Private Set-Intersection with Secure Two-Party Computation. In *SCN 18: 11th International Conference on Security in Communication Networks (Lecture Notes in Computer Science, Vol. 11035)*, Dario Catalano and Roberto De Prisco (Eds.). Springer, Heidelberg, Germany, Amalfi, Italy, 464–482. https://doi.org/10.1007/978-3-319-98113-0_25
- [15] Henry Corrigan-Gibbs and Dan Boneh. 2017. Prio: Private, Robust, and Scalable Computation of Aggregate Statistics. In *Proceedings of the 14th USENIX Conference on Networked Systems Design and Implementation (Boston, MA, USA) (NSDI'17)*. USENIX Association, USA, 259–282.
- [16] Ronald Cramer, Ivan Damgård, and Yuval Ishai. 2005. Share Conversion, Pseudorandom Secret-Sharing and Applications to Secure Computation. In *TCC 2005: 2nd Theory of Cryptography Conference (Lecture Notes in Computer Science, Vol. 3378)*, Joe Kilian (Ed.). Springer, Heidelberg, Germany, Cambridge, MA, USA, 342–362. https://doi.org/10.1007/978-3-540-30576-7_19
- [17] Morten Dahl, Jason Mancuso, Yann Dupis, Ben Decoste, Morgan Giraud, Ian Livingstone, Justin Patriquin, and Gavin Uhma. 2018. Private Machine Learning in TensorFlow using Secure Computation. *CoRR* abs/1810.08130 (2018), 1–6. arXiv:1810.08130 <http://arxiv.org/abs/1810.08130>
- [18] Dalek-Cryptography. 2020. Dalek library for elliptic curve cryptography. GitHub. <https://github.com/dalek-cryptography/curve25519-dalek>.
- [19] Ivan Damgård and Yuval Ishai. 2006. Scalable Secure Multiparty Computation. In *Advances in Cryptology – CRYPTO 2006 (Lecture Notes in Computer Science, Vol. 4117)*, Cynthia Dwork (Ed.). Springer, Heidelberg, Germany, Santa Barbara, CA, USA, 501–520. https://doi.org/10.1007/11818175_30
- [20] Hannah Davis, Christopher Patton, Mike Rosulek, and Philipp Schoppmann. 2023. Verifiable Distributed Aggregation Functions. *Proceedings on Privacy Enhancing Technologies* 2023, 4 (July 2023), 1–20.
- [21] Daniel Demmler, Thomas Schneider, and Michael Zohner. 2015. ABY - A Framework for Efficient Mixed-Protocol Secure Two-Party Computation. In *ISOC Network and Distributed System Security Symposium – NDSS 2015*. The Internet Society, San Diego, CA, USA, 1–15.
- [22] Whitfield Diffie and Martin E. Hellman. 1976. New Directions in Cryptography. *IEEE Transactions on Information Theory* 22, 6 (1976), 644–654.
- [23] Yevgeniy Dodis and Aleksandr Yampolskiy. 2005. A Verifiable Random Function with Short Proofs and Keys. In *PKC 2005: 8th International Workshop on Theory and Practice in Public Key Cryptography (Lecture Notes in Computer Science, Vol. 3386)*, Serge Vaudenay (Ed.). Springer, Heidelberg, Germany, Les Diablerets, Switzerland, 416–431. https://doi.org/10.1007/978-3-540-30580-4_28
- [24] Changyu Dong, Liqun Chen, and Zikai Wen. 2013. When private set intersection meets big data: an efficient and scalable protocol. In *ACM CCS 2013: 20th Conference on Computer and Communications Security*, Ahmad-Reza Sadeghi, Virgil D. Gligor, and Moti Yung (Eds.). ACM Press, Berlin, Germany, 789–800. <https://doi.org/10.1145/2508859.2516701>
- [25] Thai Duong, Duong Hieu Phan, and Ni Trieu. 2020. Catalic: Delegated PSI Cardinality with Applications to Contact Tracing. In *Advances in Cryptology – ASIACRYPT 2020, Part III (Lecture Notes in Computer Science, Vol. 12493)*, Shihoo Moriai and Huaxiong Wang (Eds.). Springer, Heidelberg, Germany, Daejeon, South Korea, 870–899. https://doi.org/10.1007/978-3-030-64840-4_29
- [26] Thanos Giannopoulos and Dimitris Mouris. 2018. Privacy preserving medical data analytics using secure multi party computation. An end-to-end use case. Master’s thesis. National and Kapodistrian University of Athens.
- [27] Oded Goldreich, Silvio Micali, and Avi Wigderson. 1987. How to Play any Mental Game or A Completeness Theorem for Protocols with Honest Majority. In *19th Annual ACM Symposium on Theory of Computing*, Alfred Aho (Ed.). ACM Press, New York City, NY, USA, 218–229. <https://doi.org/10.1145/28395.28420>
- [28] Mike Hamburg et al. 2020. Ristretto. <https://ristretto.group>.
- [29] Marcella Hastings, Brett Hemenway, Daniel Noble, and Steve Zdancewic. 2019. SoK: General Purpose Compilers for Secure Multi-Party Computation. In *2019 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, San Francisco, CA, USA, 1220–1237. <https://doi.org/10.1109/SP.2019.00028>
- [30] Robert Helmer, Anthony Miyaguchi, and Eric Rescorla. 2018. Testing Privacy-Preserving Telemetry with Prio. <https://hacks.mozilla.org/2018/10/testing-privacy-preserving-telemetry-with-prio>.
- [31] Mihaela Ion, Ben Kreuter, Ahmet Erhan Nergiz, Sarvar Patel, Shobhit Saxena, Karn Seth, Mariana Raykova, David Shanahan, and Moti Yung. 2020. On Deploying Secure Computing: Private Intersection-Sum-with-Cardinality. In *EuroS&P*. IEEE, Genoa, Italy, 370–389.
- [32] Marcel Keller. 2020. MP-SPDZ: A Versatile Framework for Multi-Party Computation. In *ACM CCS 2020: 27th Conference on Computer and Communications Security*, Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna (Eds.). ACM Press, Virtual Event, USA, 1575–1590. <https://doi.org/10.1145/3372297.3417872>
- [33] Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, and Laurens van der Maaten. 2021. CrypTen: Secure Multi-Party Computation Meets Machine Learning. *Advances in Neural Information Processing Systems* 34 (2021).
- [34] Vladimir Kolesnikov, Ranjit Kumaresan, Mike Rosulek, and Ni Trieu. 2016. Efficient Batched Oblivious PRF with Applications to Private Set Intersection. In *ACM CCS 2016: 23rd Conference on Computer and Communications Security*, Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi (Eds.). ACM Press, Vienna, Austria, 818–829. <https://doi.org/10.1145/2976749.2978381>
- [35] Tancrede Lepoint, Sarvar Patel, Mariana Raykova, Karn Seth, and Ni Trieu. 2021. Private Join and Compute from PIR with Default. In *Advances in Cryptology – ASIACRYPT 2021, Part II (Lecture Notes in Computer Science, Vol. 13091)*, Mehdi Tibouchi and Huaxiong Wang (Eds.). Springer, Heidelberg, Germany, Singapore, 605–634. https://doi.org/10.1007/978-3-030-92075-3_21
- [36] Yehuda Lindell and Benny Pinkas. 2009. A Proof of Security of Yao’s Protocol for Two-Party Computation. *Journal of Cryptology* 22, 2 (April 2009), 161–188. <https://doi.org/10.1007/s00145-008-9036-8>
- [37] Catherine Meadows. 1986. A More Efficient Cryptographic Matchmaking Protocol for Use in the Absence of a Continuously Available Third Party. In *1986 IEEE Symposium on Security and Privacy*. IEEE, Oakland, CA, USA, 134–134.

- <https://doi.org/10.1109/SP.1986.10022>
- [38] Peihan Miao, Sarvar Patel, Mariana Raykova, Karn Seth, and Moti Yung. 2020. Two-Sided Malicious Security for Private Intersection-Sum with Cardinality. In *Advances in Cryptology – CRYPTO 2020, Part III (Lecture Notes in Computer Science, Vol. 12172)*, Daniele Micciancio and Thomas Ristenpart (Eds.). Springer, Heidelberg, Germany, Santa Barbara, CA, USA, 3–33. https://doi.org/10.1007/978-3-030-56877-1_1
- [39] Payman Mohassel, Peter Rindal, and Mike Rosulek. 2020. Fast Database Joins and PSI for Secret Shared Data. In *ACM CCS 2020: 27th Conference on Computer and Communications Security*, Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna (Eds.). ACM Press, Virtual Event, USA, 1271–1287. <https://doi.org/10.1145/3372297.3423358>
- [40] Dimitris Mouris, Pratik Sarkar, and Nektarios Georgios Tsoutsos. 2023. PLASMA: Private, Lightweight Aggregated Statistics against Malicious Adversaries. Cryptology ePrint Archive, Report 2023/080. <https://eprint.iacr.org/2023/080>.
- [41] Dimitris Mouris and Nektarios Georgios Tsoutsos. 2021. Masquerade: Verifiable Multi-Party Aggregation with Secure Multiplicative Commitments. Cryptology ePrint Archive, Report 2021/1370. <https://eprint.iacr.org/2021/1370>.
- [42] Mahnush Movahedi, Benjamin M. Case, James Honaker, Andrew Knox, Li Li, Yiming Paul Li, Sanjay Saravanan, Shubho Sengupta, and Erik Taubeneck. 2021. Privacy-Preserving Randomized Controlled Trials: A Protocol for Industry Scale Deployment. In *Proceedings of the 2021 on Cloud Computing Security Workshop (Virtual Event, Republic of Korea) (CCSW '21)*. Association for Computing Machinery, New York, NY, USA, 59–69. <https://doi.org/10.1145/3474123.3486764>
- [43] Muhammad Naveed, Seny Kamara, and Charles V. Wright. 2015. Inference Attacks on Property-Preserving Encrypted Databases. In *ACM CCS 2015: 22nd Conference on Computer and Communications Security*, Indrajit Ray, Ninghui Li, and Christopher Kruegel (Eds.). ACM Press, Denver, CO, USA, 644–655. <https://doi.org/10.1145/2810103.2813651>
- [44] Pascal Paillier. 1999. Public-Key Cryptosystems Based on Composite Degree Residuosity Classes. In *Advances in Cryptology – EUROCRYPT'99 (Lecture Notes in Computer Science, Vol. 1592)*, Jacques Stern (Ed.). Springer, Heidelberg, Germany, Prague, Czech Republic, 223–238. https://doi.org/10.1007/3-540-48910-X_16
- [45] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. 2019. SpOT-Light: Lightweight Private Set Intersection from Sparse OT Extension. In *Advances in Cryptology – CRYPTO 2019, Part III (Lecture Notes in Computer Science, Vol. 11694)*, Alexandra Boldyreva and Daniele Micciancio (Eds.). Springer, Heidelberg, Germany, Santa Barbara, CA, USA, 401–431. https://doi.org/10.1007/978-3-030-26954-8_13
- [46] Benny Pinkas, Thomas Schneider, Oleksandr Tkachenko, and Avishay Yanai. 2019. Efficient Circuit-Based PSI with Linear Communication. In *Advances in Cryptology – EUROCRYPT 2019, Part III (Lecture Notes in Computer Science, Vol. 11478)*, Yuval Ishai and Vincent Rijmen (Eds.). Springer, Heidelberg, Germany, Darmstadt, Germany, 122–153. https://doi.org/10.1007/978-3-030-17659-4_5
- [47] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. 2018. Efficient Circuit-Based PSI via Cuckoo Hashing. In *Advances in Cryptology – EUROCRYPT 2018, Part III (Lecture Notes in Computer Science, Vol. 10822)*, Jesper Buus Nielsen and Vincent Rijmen (Eds.). Springer, Heidelberg, Germany, Tel Aviv, Israel, 125–157. https://doi.org/10.1007/978-3-319-78372-7_5
- [48] Benny Pinkas, Thomas Schneider, and Michael Zohner. 2014. Faster Private Set Intersection Based on OT Extension. In *USENIX Security 2014: 23rd USENIX Security Symposium*, Kevin Fu and Jaeyeon Jung (Eds.). USENIX Association, San Diego, CA, USA, 797–812.
- [49] Peter Rindal and Philipp Schoppmann. 2021. VOLE-PSI: Fast OPRF and Circuit-PSI from Vector-OLE. In *Advances in Cryptology – EUROCRYPT 2021, Part II (Lecture Notes in Computer Science, Vol. 12697)*, Anne Canteaut and Francois-Xavier Standaert (Eds.). Springer, Heidelberg, Germany, Zagreb, Croatia, 901–930. https://doi.org/10.1007/978-3-030-77886-6_31
- [50] Mike Rosulek and Ni Trieu. 2021. Compact and Malicious Private Set Intersection for Small Sets. In *ACM CCS 2021: 28th Conference on Computer and Communications Security*, Giovanni Vigna and Elaine Shi (Eds.). ACM Press, Virtual Event, Republic of Korea, 1166–1181. <https://doi.org/10.1145/3460120.3484778>
- [51] Adi Shamir. 1979. How to Share a Secret. *Communications of the Association for Computing Machinery* 22, 11 (Nov. 1979), 612–613.
- [52] Silvia Vermicelli, Livio Cricelli, and Michele Grimaldi. 2021. How can crowdsourcing help tackle the COVID-19 pandemic? An explorative overview of innovative collaborative practices. *R&D Management* 51, 2 (2021), 183–194.
- [53] Xiao Wang, Alex J. Malozemoff, and Jonathan Katz. 2016. EMP-toolkit: Efficient MultiParty computation toolkit. <https://github.com/emp-toolkit>.
- [54] Andrew Chi-Chih Yao. 1982. Protocols for Secure Computations (Extended Abstract). In *23rd Annual Symposium on Foundations of Computer Science*. IEEE Computer Society Press, Chicago, Illinois, 160–164. <https://doi.org/10.1109/SFCS.1982.38>
- [55] Sen Yuan, Milan Shen, Ilya Mironov, and Anderson C. A. Nascimento. 2021. Practical, Label Private Deep Learning Training based on Secure Multiparty Computation and Differential Privacy. *IACR Cryptol. ePrint Arch.* 1 (2021), 835.

APPENDIX

A ADDITIONAL DEFINITIONS

Below, we provide Algs. 1 and 2 for Defs. 3 and 4. For simplicity, our DPMC protocol in Fig. 4 uses single keys. Alg. 1 computes the join as outlined in Def. 3 for multiple keys but it can be easily adjusted for single keys. Alg. 2 computes the join as outlined in Def. 4.

Algorithm 1 Join for DPMC (Fig. 4 and Def. 3).

Match on: $\{hc_{i,i,j}\}_{j \in [m_i], i \in [T], i \in [T]}$ and $\{(h_{C,i,j})_{j \in [m_C], i \in [m_C]}\}_{i \in [m_C]}$

- 1: $\mathcal{J} := \emptyset$ ▷ Initialize join.
- 2: **For** $i \in [m_C]$, $\hat{i} \in [T]$: ▷ Perform the join.
- 3: **For** $j \in [m_{C,i}]$: ▷ Set of matched indices.
- 4: $S_{i,j,\hat{i}} := \{i' \in [m_i] \mid \exists j' \in [m_{i,i'}] \text{ s.t. } hc_{i,i',j'} = hc_{i,i,j}\}$
- 5: **If** $\bigcup_j S_{i,j,\hat{i}} \neq \emptyset$: ▷ If a match was found.
- 6: $j_{i,\hat{i}} := \min(j \in [m_{C,i}] \text{ s.t. } S_{i,j,\hat{i}} \neq \emptyset)$
- 7: Pick $i' \in S_{i,j_{i,\hat{i}},\hat{i}}$ ▷ i' is unique for each i .
- 8: Add (i', \hat{i}) to \mathcal{J} .

Algorithm 2 Join for D_s PMP (Fig. 6 and Def. 4).

Match on: $\{h_{i,j}\}_{i \in [M], j \in [m_i]}$ and $\{h_{C,i,j}\}_{i \in [m_C], j \in [m_{C,i}]}$

- 1: $\mathcal{J} := \emptyset$ ▷ Initialize join.
- 2: **For** $i \in [m_C]$: ▷ Perform the join.
- 3: **For** $j \in [m_{C,i}]$: ▷ For each column.
- 4: $S_{i,j} := \{i' \in [M] \mid \exists j' \in [m_{i'}] \text{ s.t. } h_{i',j'} = h_{C,i,j}\}$
- 5: $t_i := 1$ ▷ Keep track of number of matches for row i .
- 6: $S_T := \emptyset$
- 7: **For** $j \in [m_{C,i}]$: ▷ For each column.
- 8: **If** $\bigcup_{j \in [m_{C,i}]} S_{i,j} \setminus S_T \neq \emptyset$ and $t_i < T$: ▷ If a match was found.
- 9: $i' \stackrel{R}{\leftarrow} S_{i,j} \setminus S_T$
- 10: $S_T := S_T \cup \{i'\}$
- 11: $t_i := t_i + 1$ ▷ t_i matches for row i .
- 12: Add (i', t_i) to \mathcal{J} .

DEFINITION 14 (DDH ASSUMPTION). [22] Let $\mathbb{G}(\kappa)$ be a group parameterized by security parameter κ and g be a generator. We say that the Decisional Diffie–Hellman (DDH) assumption holds in group $\mathbb{G}(\kappa)$ if for every ppt adversary \mathcal{A} :

$$|\Pr[\mathcal{A}(g, g^a, g^b, g^{ab}) = 1] - \Pr[\mathcal{A}(g, g^a, g^b, g^c) = 1]| \leq \text{negl},$$

where the probability is taken over $a \stackrel{R}{\leftarrow} \mathbb{Z}_q$, $b \stackrel{R}{\leftarrow} \mathbb{Z}_q$, $c \stackrel{R}{\leftarrow} \mathbb{Z}_q$ and the random coins of \mathcal{A} .

DEFINITION 15 (PSEUDORANDOM GENERATOR). We call a deterministic polynomial time algorithm PRG a pseudorandom generator if for any ppt adversary \mathcal{A} ,

$$|\Pr[\mathcal{A}(x) = 1] - \Pr[\mathcal{A}(u) = 1]| \leq \text{negl},$$

where $\ell > \kappa$, $u \stackrel{R}{\leftarrow} \{0, 1\}^\ell$, seed $\stackrel{R}{\leftarrow} \{0, 1\}^\kappa$ and $x = \text{PRG}(\text{seed})$.

DEFINITION 16 (RANDOM ORACLE). [3] A random oracle RO is a family of functions that maps an input from $\{0, 1\}^*$ to an ℓ -bit image $\{0, 1\}^\ell$ s.t. each output is selected uniformly and independently.

DEFINITION 17 (SYMMETRIC KEY ENCRYPTION). A symmetric encryption scheme parameterized with security parameter κ is a triplet of algorithms (SKE.KG, SKE.Enc, SKE.Dec) with the following syntax.

- $\text{SKE.KG}(1^\kappa)$: On input 1^κ output secret key sk .
- $\text{SKE.Enc}(sk, x)$: On input (sk, x) , SKE.Enc outputs a ciphertext ct .
- $\text{SKE.Dec}(sk, ct)$: On input (sk, ct) , SKE.Dec outputs a message x .

For correctness, we ask that for any message $x \in \{0, 1\}^*$,

$$\Pr_{sk \leftarrow \text{SKE.KG}(1^\kappa)} [\text{SKE.Dec}(sk, \text{SKE.Enc}(sk, x)) = x] \geq 1 - \text{negl}.$$

DEFINITION 18 (PUBLIC KEY ENCRYPTION). A public encryption scheme parameterized with security parameter κ is a triplet of algorithms $(\text{PKE.KG}, \text{PKE.Enc}, \text{PKE.Dec})$ with the following syntax:

- $\text{PKE.KG}(1^\kappa)$: On input 1^κ output a key pair (pk, sk) .
- $\text{PKE.Enc}(pk, x)$: On input (pk, x) , PKE.Enc outputs a ciphertext ct .
- $\text{PKE.Dec}(sk, ct)$: On input (sk, ct) , PKE.Dec outputs a message x .

For correctness, we ask that for any message $x \in \{0, 1\}^*$,

$$\Pr_{(pk, sk) \leftarrow \text{PKE.KG}(1^\kappa)} [\text{PKE.Dec}(sk, \text{PKE.Enc}(pk, x)) = x] \geq 1 - \text{negl}.$$

B RERANDOMIZABLE ENCRYPTED OPRF (EO)

B.1 EO Definition

In Def. 11, we introduce a new construction called rerandomizable encrypted OPRF (EO) that allows two parties to encrypt, mask, and shuffle their data.

DEFINITION 19 (PSEUDORANDOMNESS OF THE EVALUATION). We say that the evaluation is pseudorandom if for any ppt adversary \mathcal{A} with query access to $\mathcal{O}_{\text{Eval}(sk, \cdot)}$ ($\mathcal{O}_u(\cdot)$),

$$|\Pr[\mathcal{A}^{\mathcal{O}_{\text{Eval}(sk, \cdot)}}(pk, pf) = 1] - \Pr[\mathcal{A}^{\mathcal{O}_u(\cdot)}(pk, pf) = 1]| \leq \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$, $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$, and for $x \in \{0, 1\}^\kappa$, \mathcal{O}_u outputs a uniform y whereas $\mathcal{O}_{\text{Eval}(sk, \cdot)}$ outputs $y = \text{Eval}(sk, x)$.

A stronger definition of pseudorandomness of the evaluation is malicious pseudorandomness of the oblivious evaluation. We add the definition for completeness even though our construction only satisfies the pseudorandomness of the evaluation.

DEFINITION 20 (MALICIOUS PSEUDORANDOMNESS OF THE OBLIVIOUS EVALUATION). We say that the oblivious evaluation is pseudorandom if for any ppt adversary \mathcal{A} with query access to $\mathcal{O}_{\text{Dec}(sk, \text{OEval}(ek, \cdot))}$ ($\mathcal{O}_u(\cdot)$),

$$|\Pr[\mathcal{A}^{\mathcal{O}_{\text{Dec}(sk, \text{OEval}(ek, \cdot))}}(pk, pf) = 1] - \Pr[\mathcal{A}^{\mathcal{O}_u(\cdot)}(pk, pf) = 1]| \leq \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$, $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$, and for $ct \leftarrow \mathcal{A}^{\mathcal{O}}(pk, pf)$ with $\text{Dec}(sk, \text{OEval}(ek, ct)) \neq \perp$, \mathcal{O}_u outputs a uniform y whereas $\mathcal{O}_{\text{Dec}(sk, \text{OEval}(ek, \cdot))}$ outputs $y = \text{Dec}(sk, \text{OEval}(ek, ct))$.

DEFINITION 21 (CIPHERTEXT INDISTINGUISHABILITY FOR EVALUATION KEY (ek) OWNER). We call EO ciphertext indistinguishable for the evaluation key owner if for any ppt algorithm \mathcal{A} ,

$$|\Pr[\mathcal{A}(pk, ct_0) = 1] - \Pr[\mathcal{A}(pk, ct_1) = 1]| \leq \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$. In the adaptive malicious setting $(m_0, m_1, pf) \leftarrow \mathcal{A}(pk)$ whereas in the semi-honest setting $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$ and $(m_0, m_1) \leftarrow \mathcal{A}(pk, pf, ek)$. $\forall i \in \{0, 1\} : ct_i \leftarrow \text{Enc}(pk, pf, m_i)$.

DEFINITION 22 (CIPHERTEXT INDISTINGUISHABILITY FOR SECRET KEY OWNER). We call EO ciphertext indistinguishable for the secret key owner if for any ppt algorithm \mathcal{A} ,

$$|\Pr[\mathcal{A}(pk, ct_0) = 1] - \Pr[\mathcal{A}(pk, ct_1) = 1]| \leq \text{negl},$$

where $(pf, ek) \leftarrow \text{KG}(1^\kappa)$. In the adaptive malicious setting $(m_0, m_1, pk) \leftarrow \mathcal{A}(pf)$ whereas in the semi-honest setting $(pk, sk) \leftarrow \text{KG}(1^\kappa)$ and $(m_0, m_1) \leftarrow \mathcal{A}(pk, pf, sk)$. $\forall i \in \{0, 1\} : ct_i \leftarrow \text{Enc}(pk, pf, m_i)$.

DEFINITION 23 (RERANDOMIZED CIPHERTEXT INDISTINGUISHABILITY). We call EO rerandomized ciphertext indistinguishable if for any ppt algorithm \mathcal{A} ,

$$|\Pr[\mathcal{A}(pk, ct_0) = 1] - \Pr[\mathcal{A}(pk, ct_1) = 1]| \leq \text{negl},$$

$(x, pk, pf) \leftarrow \mathcal{A}(1^\kappa)$, $ct_0 \leftarrow \text{Rnd}(pk, pf, \text{Enc}(pk, pf, x))$ and $ct_1 \leftarrow \text{Enc}(pk, pf, x)$.

DEFINITION 24 (CIPHERTEXT WELL-FORMEDNESS). We call an EO scheme ciphertext well-formed if for any x_0, x_1 with $\text{OEval}(ek, x_0) = \text{OEval}(ek, x_1)$

$$\Delta_s(ct_0, ct_1) \leq \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$, $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$ and Δ_s is the statistical distance.

DEFINITION 25 (EVALUATED CIPHERTEXT SIMULATABILITY). We call an EO scheme evaluated ciphertext simulatable if there exists an ppt algorithm EO.Sim such that for any x ,

$$\Delta_s(\text{ect}_0, \text{ect}_1) \leq \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$, $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$, $\text{ect}_0 \leftarrow \text{OEval}(ek, \text{Enc}(pk, pf, x))$, $\text{ect}_1 \leftarrow \text{EO.Sim}(pk, pf, sk, \text{Eval}(ek, x))$ and Δ_s is the statistical distance.

B.2 EO Construction and Security Analysis

In this section, we instantiate our EO construction in cyclic groups and prove its security against semi-honest adversaries.

DEFINITION 26 (EO CONSTRUCTION IN CYCLIC GROUPS). Let g be a generator of a cyclic group \mathbb{G} with order q and $H_{\mathbb{G}}(\cdot) : \{0, 1\}^* \rightarrow \mathbb{G}$ a hash function. Then the EO collection of algorithms is constructed as follows.

- $\text{KG}(1^\kappa)$: Sample $a \xleftarrow{\mathbb{R}} \mathbb{Z}_q$ and output $(pk := g^a, sk := a)$.
- $\text{EKG}(1^\kappa)$: Sample $b \xleftarrow{\mathbb{R}} \mathbb{Z}_q$ and output $(pf := g^b, ek := b)$.
- $\text{Eval}(ek, x)$: Output $y = H_{\mathbb{G}}(x)^{ek}$.
- $\text{Enc}(pk, pf, x)$: Sample $r \xleftarrow{\mathbb{R}} \mathbb{Z}_q$ and define $ct_1 := pf^r$, $ct_2 := pk^r \cdot H_{\mathbb{G}}(x)$. If $pk \neq pf$ output ciphertext $ct := (ct_1, ct_2)$ otherwise output \perp .
- $\text{Rnd}(pk, pf, ct)$: Let $ct = (ct_1, ct_2)$. Sample $r \xleftarrow{\mathbb{R}} \mathbb{Z}_q$ and define $ct'_1 := ct_1 \cdot pf^r$, $ct'_2 := ct_2 \cdot pk^r$ and output ciphertext $ct' := (ct'_1, ct'_2)$.
- $\text{OEval}(ek, ct)$: Let $ct = (ct_1, ct_2)$. Define $ect_2 := ct_2^{ek}$ and output $ect := (ct_1, ect_2)$.
- $\text{Dec}(sk, ect)$: Let $ect = (ect_1, ect_2)$. Output $y := ect_2 / ect_1^{sk}$.

For correctness, we ask that for any $x \in \{0, 1\}^*$,

$$\Pr[\text{Dec}(sk, \text{OEval}(ek, \text{Rnd}(pk, pf, \text{Enc}(pk, pf, x)))) = \text{Eval}(ek, x)] \geq 1 - \text{negl},$$

where $(pk, sk) \leftarrow \text{KG}(1^\kappa)$ and $(pf, ek) \leftarrow \text{EKG}(1^\kappa)$.

LEMMA 27. *Def. 26 defines a correct EO scheme.*

PROOF. Let $(g^a, a) \leftarrow \text{KG}(1^\kappa)$ and $(g^b, b) \leftarrow \text{EKG}(1^\kappa)$. The correctness of the EO construction is satisfied as shown below:

$$\begin{aligned} \text{Dec}(a, \text{OEval}(b, \text{Rnd}(g^a, g^b, \text{Enc}(g^a, g^b, x)))) &= \text{Eval}(b, x) \Leftrightarrow \\ \text{Dec}(a, \text{OEval}(b, \text{Rnd}(g^a, g^b, (g^{br}, g^{ar} \cdot H_{\mathbb{G}}(x)))))) &= H_{\mathbb{G}}(x)^b \Leftrightarrow \\ \text{Dec}(a, \text{OEval}(b, (g^{br} \cdot g^{br'}, g^{ar} \cdot g^{ar'} \cdot H_{\mathbb{G}}(x)))) &= H_{\mathbb{G}}(x)^b \Leftrightarrow \\ \text{Dec}(a, \text{OEval}(b, (g^{b(r+r')}, g^{a(r+r')} \cdot H_{\mathbb{G}}(x)))) &= H_{\mathbb{G}}(x)^b \Leftrightarrow \\ \text{Dec}(a, (g^{b(r+r')}, (g^{a(r+r')} \cdot H_{\mathbb{G}}(x))^b)) &= H_{\mathbb{G}}(x)^b \Leftrightarrow \\ \text{Dec}(a, (g^{b(r+r')}, g^{ab(r+r')} \cdot H_{\mathbb{G}}(x)^b)) &= H_{\mathbb{G}}(x)^b \Leftrightarrow \\ g^{ab(r+r')} \cdot H_{\mathbb{G}}(x)^b / (g^{b(r+r')})^a &= H_{\mathbb{G}}(x)^b. \end{aligned}$$

□

The construction is secure against semi-honest adversaries under the DDH assumption. The bottleneck that prevents malicious security is the OPRF $H(x)^k$. This OPRF only provides semi-honest security since a malicious delegator might send an arbitrary group element X instead of $H(x)$. In that case, it does not result in an OPRF since it satisfies linear relations, e.g., $X^k \cdot Y^k = (X \cdot Y)^k$.

We have outlined what is needed from the EO for malicious security in Defs. 20-22. The main bottleneck for our $H(x)^k$ based construction is Def. 20 (Defs. 21 and 22 seem to hold when making stronger assumptions than DDH). Other PRF candidates seem significantly less efficient (i.e., lowMC) or require stronger assumptions (e.g., Dodis-Yampolskiy PRF [23]). In Appendix C, we show that our EO primitive is compatible with the MPC shuffle protocol of [39] by relying on the EO rerandomization procedure.

LEMMA 28. *Def. 26 satisfies pseudorandomness of the evaluation under the DDH assumption in the Random Oracle Model.*

PROOF. We use a sequence of hybrids in which we replace step by step (based on the order of random oracle queries) $\text{Eval}(ek, x)$ with a uniform group element. If there is a distinguisher against the pseudorandomness of Eval with probability ϵ then there is a distinguisher against at least two consecutive intermediate hybrids with probability ϵ/Q , where Q is the maximum between the amount of random oracle and Eval oracle queries. Given such a distinguisher, we build a distinguisher against DDH as follows. The DDH distinguisher receives challenge A, B, C and sets $\text{pf} := A$. Once the random oracle query is made that differentiates the two hybrids (let that be the i^* th query), it programs $H_{\mathbb{G}}(x) := B$. For all following queries $i > i^*$ program $H_{\mathbb{G}}(x) := g^{r_i}$, where $r_i \xleftarrow{R} \mathbb{Z}_q$. When a query for x to the Eval oracle is made, query x to the random oracle if it has not been made yet. If x matches the query i^* , respond with C . If x corresponds to a query $i < i^*$, respond with a uniform group element. Otherwise respond with B^{r_i} .

If $A = g^a, B = g^b, C = g^c$ then the DDH distinguisher simulates the first of the two hybrids. In case of uniform A, B, C it simulates the second of the two hybrids where the output of the Eval oracle that corresponds to the i^* th message is uniform.

Since Q is polynomial and the distinguishing probability against DDH is negligible, the probability to break the pseudorandomness of Eval is also negligible. □

LEMMA 29. *Def. 26 is ciphertext indistinguishable for the evaluation key owner in the semi-honest setting under the DDH assumption.*

PROOF. We use three hybrids, the first hybrid uses x_0 for the challenge ciphertext. In the second hybrid, the ciphertext is independent of the message. The third hybrid uses x_1 for the challenge ciphertext. We show now that these three hybrids cannot be distinguished based on the DDH assumption.

We build a DDH distinguisher for hybrid one and two (two and three) as follows. It receives DDH challenge A, B, C and samples $(\text{pk}, \text{ek}) \leftarrow \text{EKG}(1^\kappa)$. It defines $\text{pk} := A$ and sends $(\text{pk}, \text{ek}, \text{pf})$ to the distinguisher against the ciphertext indistinguishability. It receives x_0 and x_1 . Return challenge ciphertext $\text{ct}_1 := B^{\text{ek}}, \text{ct}_2 := C \cdot x_0$ ($\text{ct}_2 := C \cdot x_1$). Output the output of the ciphertext indistinguishability distinguisher.

If $A = g^a, B = g^b, C = g^c$ then the challenge ciphertext follows the output distribution of Enc for x_0 as in the first hybrid (and m_1 in the third hybrid). Otherwise, the challenge ciphertext is independent of the message as in the second hybrid. □

LEMMA 30. *Def. 26 is ciphertext indistinguishable for the secret key owner in the semi-honest setting under the DDH assumption for prime groups (every element is a generator).*

PROOF. We use three hybrids, the first hybrid uses x_0 for the challenge ciphertext. In the second hybrid, the ciphertext is independent of the message. The third hybrid uses x_1 for the challenge ciphertext. We show now that these three hybrids cannot be distinguished based on the DDH assumption.

We build a DDH distinguisher for hybrid one and two (two and three) as follows. It receives DDH challenge A, B, C and samples $(\text{pk}, \text{sk}) \leftarrow \text{KG}(1^\kappa)$. It defines $\text{pf} := A$ and sends $(\text{pk}, \text{sk}, \text{pf})$ to the distinguisher against the ciphertext indistinguishability. It receives x_0 and x_1 . Return challenge ciphertext $\text{ct}_1 := C, \text{ct}_2 := B^{\text{sk}} \cdot x_0$ ($\text{ct}_2 := B^{\text{sk}} \cdot x_1$). Output the output of the ciphertext indistinguishability distinguisher.

If $A = g^a, B = g^b, C = g^c$ then the challenge ciphertext follows the output distribution of Enc for x_0 (x_1) as in the first hybrid (third hybrid). Otherwise, the challenge ciphertext is independent of the message as in the second hybrid as long as B is a generator of the group and thus B^{sk} is uniform for a uniform B . □

LEMMA 31. *Def. 26 is statistically randomized ciphertext indistinguishable.*

PROOF. Let $\text{ct} := (g^{br}, g^{ar} \cdot H_{\mathbb{G}}(x))$ be an encryption of x for some random $r \in \mathbb{Z}_q$. Then the randomized ciphertext $\text{Rnd}(g^a, g^b, \text{ct})$ is defined as $(g^{br} \cdot g^{br'}, g^{ar} \cdot g^{ar'} \cdot H_{\mathbb{G}}(x)) = (g^{b(r+r')}, g^{a(r+r')} \cdot H_{\mathbb{G}}(x))$ for random $r' \in \mathbb{Z}_q$. Since both r and r' are random elements in \mathbb{Z}_q , $r + r'$ is also a random element in \mathbb{Z}_q and the ciphertext is statistically randomized ciphertext indistinguishable. □

LEMMA 32. *Let sk and q be coprime. Then Def. 26 is ciphertext well formed.*

PROOF. Ciphertext well-formedness demands that messages that result in the same PRF evaluation have an identical ciphertext distribution. In the construction of Def. 26 the ciphertext only depends on

$H_{\mathbb{G}}(x)$ and the output of Eval is $H_{\mathbb{G}}(x)^{ek}$. Now, let there be x_0 and x_1 with $H_{\mathbb{G}}(x_0)^{ek} = H_{\mathbb{G}}(x_1)^{ek}$ and let for $b \in \{0, 1\}$, $H_{\mathbb{G}}(x_b) = g^{r^b}$. Then $(r - r') \cdot ek = 0 \pmod q$ and therefore $(r - r') = 0$ such that $H_{\mathbb{G}}(x_0) = H_{\mathbb{G}}(x_1)$ and the ciphertexts have the same distribution or ek would divide the group order q and therefore not be coprime. \square

LEMMA 33. *Def. 26 is evaluated ciphertext simulatable.*

PROOF. EO.Sim takes as input $pk = g^a$, $pf = g^b$, $sk = a$ and $y = H_{\mathbb{G}}(x)^b$. It outputs $ect = (ect_0, ect_1)$ where $ect_0 = pf^r$, $ect_1 = pf^{ar}$. y . This is identically distributed as $ect = \text{OEval}(ek, \text{Enc}(pk, pf, x)) = (g^{r'b}, g^{r'ab} \cdot H_{\mathbb{G}}(x))$. \square

THEOREM 34. *Def. 26 is a secure and correct EO scheme. More precisely, it is correct, satisfies pseudorandomness of the evaluation and ciphertext well-formedness, evaluated ciphertext simulatability, is randomized ciphertext indistinguishable as well as ciphertext indistinguishable for the evaluation and secret key owner. The latter two are semi-honest secure under the DDH assumption.*

PROOF. Follows from Lemma 27, 28, 29, 30, 31, 32, and 33. \square

C THREE-PARTY SECURE SHUFFLING FOR D_S PMC

C.1 Ideal Shuffle Functionality

The ideal shuffle functionality from Fig. 12 gets inputs from parties C and D secret shares and generates fresh shuffled shares and sends them back to parties S and D . Additionally, $\mathcal{F}_{\text{Shuffle}}$ gets multiple EO ciphertexts from C , generates fresh shuffled ciphertexts, and sends them back to C . Parties P_1 to P_T do not participate in the protocol but do have information about the encrypted and secret shared information and might be corrupted.

C.2 Shuffle Protocol

We define a permutation of size m_C as an injective function $\pi : [N] \rightarrow [N]$. We denote as π_{AB} a permutation generated from party A and sent to B . Fig. 13 demonstrates the honest majority shuffling protocol utilized by D_S PMC. Our shuffling protocol performs two iterations of a permutation network and reshapes C 's and D 's inputs (sh_C and sh_D , respectively). Parties C and D have T sh_C and sh_D vectors (indicated as $sh_{C,t}$, $sh_{D,t}$ for $t \in [T]$), each of which has m_t elements. Additionally, the shuffling protocol reshapes EO.ct to prevent leakage of honest parties' data in the presence of an adversary that has corrupted D and multiple parties P .

The first iteration of the permutation network is demonstrated in steps 1-3 in Fig. 13 and reshapes sh_C , sh_D to \overline{SH}_C , \overline{SH}_S and EO.ct to $\overline{EO.CT}$. Party C generates two permutations (π_{CS} and π_{CD}) as well as two vectors of scalars (V_{CS} and V_{CD}) to rerandomize sh_C and sh_D . C locally applies the two permutations and XORs with the vectors of scalars. C then sends one permutation and one vector of scalars to each of D and S . D first permutes and XORs sh_D with V_{CD} and sends the result to S who, in turn, permutes it with π_{CS} and XORs it with V_{CS} to compute \overline{SH}_S .

In the second iteration, party S generates two more permutations (π_{SC} and π_{SD}) as well as two vectors of scalars (V_{SC} and V_{SD}) to rerandomize the outputs of the first iteration (i.e., \overline{SH}_C and \overline{SH}_S). Next, S applies both permutations on \overline{SH}_S and XORs it with both

vectors V_{SC} and V_{SD} , while parties C and D communicate to apply the same operations on \overline{SH}_C . At the end of the protocol, S gets \overline{SH}_C and D gets \overline{SH}_D such that $\overline{SH}_C \oplus \overline{SH}_D = sh_C \oplus sh_D$. Finally, party C gets $\overline{EO.CT}$, which is the blinded and rerandomized EO.ct.

Observe that the communication in the aforementioned protocol is only linear to the size of EO.CT. We can further optimize the communication by having each two parties (C with D , C with S , and S with D) pre-share some randomness and use it as a PRF key. These PRF keys can then be used to generate both the random permutations and the random vectors of scalars which will be consistent between the parties.

C.3 Security Analysis of Secure Shuffling

THEOREM 35. *Let EO be a correct Rerandomizable Encrypted OPRF scheme that satisfies statistical rerandomized ciphertext indistinguishability, ciphertext indistinguishability (for evaluation key or secret key owner) and ciphertext well-formedness. Then, the shuffling protocol in Fig. 13 realizes the ideal shuffling functionality in Fig. 12 when at most one of the parties C , D and S and any amount of the parties P_1 to P_t are corrupted and semi-honest.*

PROOF. We prove the theorem by showing that for each party, there exists a simulator that produces a view that is indistinguishable from the view of the corrupted party in the real shuffle protocol.

CLAIM 36. *Let EO be correct, satisfy statistical randomized ciphertext indistinguishability and ciphertext well-formedness. Then, there is a simulator that produces a view of Party C that is indistinguishable from the real view of Party C for any amount of corrupted parties P_1 to P_T . We emphasize that the distinguisher also receives the in and outputs to and from the ideal functionality (which is identical to the real output) of the honest parties.*

PROOF. We first show the simulator in case none of the parties P_1 to P_T is corrupted. The view of Party C can be generated from its input $\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}, sh_{C,t,i}\}_{i \in [m_t], t \in [T]}\}$, output $\overline{EO.CT}$ and the message $\pi_{SC}, V_{SC}, \overline{EO.CT}$ from Party S . Our simulator emulates these messages and otherwise follows the description of the computation of Party C .

Our simulator receives input $\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}, sh_{C,t,i}\}_{i \in [m_t], t \in [T]}\}$, $\overline{EO.CT}$ and generates Party S 's message as follows. It uses $\overline{EO.CT}$ that was part of the input and samples $\pi_{SC} \xleftarrow{R} \text{Perm}(M)$ and $V_{SC} \xleftarrow{R} \{0, 1\}^{M \cdot |v|}$.

We now show that this simulator emulates the correct distribution. Let $\overline{EO.CT} = \pi(\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}\}_{i \in [m_t], t \in [T]}) = \pi'_{SD}(\pi'_{SC}(\pi'_{CS}(\pi'_{CD}(\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}\}_{i \in [m_t], t \in [T]}\})))$, where π'_{SD} , π'_{SC} , π'_{CS} , π'_{CD} are defined as in the original protocol and π'_{SC} , π'_{CS} , π'_{CD} are part of party C 's view. Sampling $\pi'_{SD}, \pi'_{SC}, \pi'_{CS}, \pi'_{CD} \xleftarrow{R} \text{Perm}(M)$ and defining π as their composition results in the same distribution as when sampling $\pi, \pi'_{SC}, \pi'_{CS}, \pi'_{CD} \xleftarrow{R} \text{Perm}(M)$ and defining π'_{SD} such that it is consistent with the protocol specification. The former is the distribution during a real protocol execution while the later is the distribution during the simulated run where the ideal functionality samples π and the simulator samples $\pi'_{SC}, \pi'_{CS}, \pi'_{CD}, \pi$ and π'_{SD} remain hidden from the view of Party C .

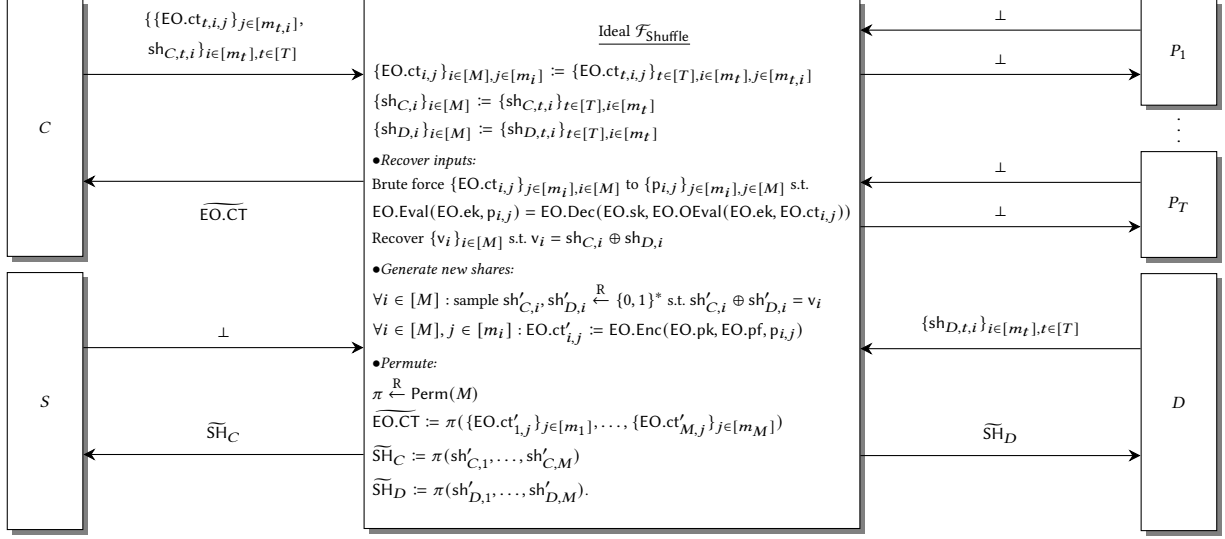


Figure 12: The figure shows the ideal $\mathcal{F}_{\text{Shuffle}}$ functionality. We define $M := \sum_{t=1}^T m_t$. We treat EO.pk and EO.pf as publicly known to all parties. Party C has access to EO.ek and Party D to EO.sk . Further, any amount of Parties P_1 to P_T can be corrupted who have access to $\{\text{EO.ct}_{t,i,j}\}_{j \in [m_t], i \in [M]}$, $\{\text{sh}_{C,t,i}\}_{i \in [m_t], t \in [T]}$ and $\{\text{sh}_{D,t,i}\}_{i \in [m_t], t \in [T]}$.

We follow this argument for the distribution of V_{SC} . There exists a unique $V \in \{0, 1\}^{M \cdot |v|}$ such that $\text{SH}_D = \text{SH}'_D \oplus V$, where SH_D denotes the original shares sent by Party D to the ideal functionality and SH'_D are the shares generated and output by the ideal functionality. The same holds for SH_C and SH'_C . Further, as specified by the protocol V can also be defined as $V := V_{SD} \oplus V_{SC} \oplus V_{CS} \oplus V_{CD}$. Here we ignore the fact that V is actually impacted by the permutations $\pi_{SD}, \pi_{SC}, \pi_{CS}, \pi_{CD}$ since it can simply be accounted for by permuting $V_{SD}, V_{SC}, V_{CS}, V_{CD}$. Both definitions of V are consistent since any two two out of two secret shares result in the same shares up to an offset vector in $\{0, 1\}^{M \cdot |v|}$. As previously sampling first $V_{SD}, V_{SC}, V_{CS}, V_{CD}$ results in the same distribution as sampling first $V, V_{SC}, V_{CS}, V_{CD}$.

The last part to show is that the output $\widetilde{\text{EO.CT}}$ of the ideal functionality is identically distributed as the Party C 's output in the real execution. From the statistical randomized ciphertext indistinguishability of EO follows that any rerandomized ciphertext for input $p_{i,j}$ is indistinguishable from a fresh encryption of $p_{i,j}$ even when given EO.ek (and EO.sk). Using a hybrid argument over all $N = \sum_{i=1}^M (m_i)$ (i.e., M total rows and each row i has m_i identifiers) distinguishing the real from the simulated view with advantage ϵ results in a ϵ/N distinguishing advantage in the randomized ciphertext indistinguishability game. Now, we show that brute forcing a $p'_{i,j}$ from a ciphertext and encrypting it is except negligible probability identically distributed as a ciphertext of $p_{i,j}$. By the correctness property it follows that except negligible probability, both ciphertexts evaluate to the same OPRF evaluation, i.e., $\text{EO.Eval}(\text{EO.ek}, p_{i,j}) = \text{EO.Eval}(\text{EO.ek}, p'_{i,j})$. Now, we can invoke

the ciphertext well-formedness which ensures that the rerandomized $\widetilde{\text{EO.CT}}$ is with overwhelming probability identically distributed as the fresh $\widetilde{\text{EO.CT}}$ generated by the ideal functionality.

In case some of the parties P_1 to P_T are corrupted we actually do not need to adapt our simulator. The difference is that when adding the views of the corrupted parties among P_1 to P_T to the view of C , Party C has access to some of the shares $\{\text{SH}_{D,t,i}\}_{i \in [m_t], t \in [T]}$. However, knowing these shares do not have impact on the distribution of the view generated by our simulator and can therefore simply added to the view. \square

CLAIM 37. *There exists a simulator that produces a view of Party D that is indistinguishable from the real view of Party D for any amount of corrupted parties P_1 to P_T .*

PROOF. We start with the case where there is no corruption among parties P_1 to P_T . Party D 's view can be generated from its input $\{\text{sh}_{D,t,i}\}_{i \in [m_t], t \in [T]}$, output $\widetilde{\text{SH}}_D$ and the messages (π_{CD}, V_{CD}) , $\widetilde{\text{SH}}_C$ from Party C and π_{SD}, V_{SD} from Party S . Therefore it suffices for our simulator to emulate these messages and generate the view from these messages according to the protocol description.

Our simulator on input $\{\text{sh}_{D,t,i}\}_{i \in [m_t], t \in [T]}$, $\widetilde{\text{SH}}_D$ samples $\pi_{CD}, \pi_{SD} \xleftarrow{R} \text{Perm}(M), V_{CD}, V_{SD} \xleftarrow{R} \{0, 1\}^{M \cdot |v|}$. $\widetilde{\text{SH}}_C$ is picked such that $\widetilde{\text{SH}}_D = \pi_{SD}(\widetilde{\text{SH}}_C) \oplus V_{SD}$. We define π as in the previous claim. As previously, sampling first $\pi_{SD}, \pi_{SC}, \pi_{CS}, \pi_{CD}$ and defining π as their composition as done in the real protocol execution results in the same distribution as when sampling π, π_{SD}, π_{CD} first and then defining and sampling π_{SC}, π_{CS} (not part of the view) such that they are consistent with the real protocol distribution. Using

We define $M := \sum_{t=1}^T m_t$.

<p>① First Shuffling (Party C)</p> <p>Input: $\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}, \text{sh}_{C,t,i}\}_{i \in [m_t], t \in [T]}$</p> <p>1: $V_{CD}, V_{CS} \xleftarrow{R} \{0, 1\}^{M \cdot v }$</p> <p>2: $\pi_{CD}, \pi_{CS} \xleftarrow{R} \text{Perm}(M)$</p> <p>3: For $t \in [T], i \in [m_t], j \in [m_{t,i}]$ ▷ Randomize</p> <p>4: $\text{EO.ct}'_{t,i,j} := \text{EO.Rnd}(\text{EO.pk}, \text{EO.pf}, \text{EO.ct}_{t,i,j})$</p> <p>5: $\text{SH}_C := (\text{sh}_{C,1,1}, \dots, \text{sh}_{C,m_T,T})$</p> <p>6: $\text{EO.CT} := (\{\text{EO.ct}'_{1,1,j}\}_j, \dots, \{\text{EO.ct}'_{T,m_T,j}\}_j)$</p> <p>7: $\widetilde{\text{SH}}_C := \pi_{CS}(\pi_{CD}(\text{SH}_C) \oplus V_{CD}) \oplus V_{CS}$ ▷ Perm. & Rand.</p> <p>8: $\widetilde{\text{EO.CT}} := \pi_{CS}(\pi_{CD}(\text{EO.CT}))$ ▷ Permute</p> <p>Send to S: $\pi_{CS}, V_{CS}, \widetilde{\text{EO.CT}}$</p> <p>Send to D: π_{CD}, V_{CD}</p> <p>Output of first shuffle: $\widetilde{\text{SH}}_C$</p>	<p>④ Second Shuffling (Party S)</p> <p>Input: $\text{SH}_S, \text{EO.CT}$</p> <p>1: $(\{\text{EO.ct}_{1,j}\}_j, \dots, \{\text{EO.ct}_{M,j}\}_j) := \widetilde{\text{EO.CT}}$</p> <p>2: $V_{SC}, V_{SD} \xleftarrow{R} \{0, 1\}^{M \cdot v }$</p> <p>3: $\pi_{SC}, \pi_{SD} \xleftarrow{R} \text{Perm}(M)$</p> <p>4: For $i \in [M], j \in [m_i]$ ▷ Randomize</p> <p>5: $\text{EO.ct}'_{i,j} := \text{EO.Rnd}(\text{EO.pk}, \text{EO.pf}, \text{EO.ct}_{i,j})$</p> <p>6: $\widetilde{\text{EO.CT}} := (\{\text{EO.ct}'_{1,j}\}_j, \dots, \{\text{EO.ct}'_{M,j}\}_j)$</p> <p>7: $\widetilde{\text{SH}}_C = \pi_{SD}(\pi_{SC}(\text{SH}_S) \oplus V_{SC}) \oplus V_{SD}$ ▷ Perm. & Rand.</p> <p>8: $\text{EO.CT} := \pi_{SD}(\pi_{SC}(\widetilde{\text{EO.CT}}))$ ▷ Permute</p> <p>Send to C: $\pi_{SC}, V_{SC}, \widetilde{\text{EO.CT}}$</p> <p>Send to D: π_{SD}, V_{SD}</p> <p>Output: $\widetilde{\text{SH}}_C$</p>
<p>② First Shuffling (Party D)</p> <p>Input: $\{\text{sh}_{D,t,i}\}_{i \in [m_t], t \in [T]}$</p> <p>Messages: V_{CD}, π_{CD}</p> <p>1: $\text{SH}_D := (\text{sh}_{D,1,1}, \dots, \text{sh}_{D,T,m_T})$</p> <p>2: $\widetilde{\text{SH}}_D := \pi_{CD}(\text{SH}_D) \oplus V_{CD}$ ▷ Permute and Randomize</p> <p>Send to S: $\widetilde{\text{SH}}_D$</p> <p>Output of first shuffle: -</p>	<p>⑤ Second Shuffling (Party C)</p> <p>Input: SH_C</p> <p>Messages: $\pi_{SC}, V_{SC}, \widetilde{\text{EO.CT}}$</p> <p>1: $\widetilde{\text{SH}}_C := \pi_{SC}(\text{SH}_C) \oplus V_{SC}$ ▷ Permute and Randomize.</p> <p>Send to D: $\widetilde{\text{SH}}_C$</p> <p>Output: $\widetilde{\text{EO.CT}}$</p>
<p>③ First Shuffling (Party S)</p> <p>Input: -</p> <p>Messages: $\pi_{CS}, V_{CS}, \widetilde{\text{SH}}_D, \widetilde{\text{EO.CT}}$</p> <p>1: $\text{SH}_S := \pi_{CS}(\widetilde{\text{SH}}_D) \oplus V_{CS}$ ▷ Permute and Randomize</p> <p>Output of first shuffle: $\text{SH}_S, \widetilde{\text{EO.CT}}$</p>	<p>⑥ Second Shuffling (Party D)</p> <p>Input: -</p> <p>Messages: $\pi_{SD}, V_{SD}, \widetilde{\text{SH}}_C$</p> <p>1: $\widetilde{\text{SH}}_D := \pi_{SD}(\widetilde{\text{SH}}_C) \oplus V_{SD}$ ▷ Permute and Randomize</p> <p>Output: $\widetilde{\text{SH}}_D$</p>

Figure 13: Three-Party Shuffling. Parties C and D get secret shares sh_C and sh_D of a vector v as inputs such that $v = \text{sh}_C \oplus \text{sh}_D$. Party C additionally inputs a Rerandomizable Encrypted OPRF ciphertext vector EO.ct of same length as sh_C and sh_D . The protocol reshares $(\text{sh}_C, \text{sh}_D)$ to $(\widetilde{\text{SH}}_C, \widetilde{\text{SH}}_D)$ and carries along EO.ct and reshares it to $\widetilde{\text{EO.CT}}$.

the same approach, we can show that V_{CD}, V_{SD} are also correctly distributed.

Similar to the previous claim, corrupting any amount of parties P_1 to P_T and adding them to the view of Party D does not impact the distribution of the view generated by the simulator. Again, we can simply add $\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}, \text{sh}_{C,t,i}\}_{i \in [m_t], t \in [T]}$ of the corrupted parties to the view generated by our simulator. \square

CLAIM 38. *Let EO satisfy statistical rerandomized ciphertext indistinguishability and ciphertext indistinguishability (for evaluation key or secret key owner). Then, there is a simulator that produces a view of Party S that is indistinguishable from the real view of Party S for any amount of corrupted parties P_1 to P_T .*

PROOF. Let $\{\{\text{EO.ct}'_{t,i,j}\}_{j \in [m_{t,i}]}, \text{sh}'_{C,t,i}, \text{sh}'_{D,t,i}\}_{i \in [m_t], t \in \mathbb{C} \subseteq [T]}$ be the views of the corrupted parties among P_1 to P_T . The view of Party S and the corrupted parties among P_1 and P_T can be generated from $\{\{\text{EO.ct}'_{t,i,j}\}_{j \in [m_{t,i}]}, \text{sh}'_{C,t,i}, \text{sh}'_{D,t,i}\}_{i \in [m_t], t \in \mathbb{C} \subseteq [T]}$, S 's output $\widetilde{\text{SH}}_C$, the messages $\pi_{CS}, V_{CS}, \widetilde{\text{EO.CT}}$ from Party D and $\widetilde{\text{SH}}_D$ from Party D .

Our simulator has inputs $\{\{\text{EO.ct}'_{t,i,j}\}_{j \in [m_{t,i}]}, \text{sh}'_{C,t,i}, \text{sh}'_{D,t,i}\}_{i \in [m_t], t \in \mathbb{C} \subseteq [T]}$ and $\widetilde{\text{SH}}_C$. It samples $\pi_{CS} \xleftarrow{R} \text{Perm}(M), V_{CS} \xleftarrow{R} \{0, 1\}^{M \cdot |v|}$ and generates $\widetilde{\text{EO.CT}}$ as encryptions of 0.

We now show that the simulator generates the correct distribution. We define π as in the previous claims. $\pi_{SD}, \pi_{SC}, \pi_{CS}, \pi_{CD}$ are part of the simulated view except π and π_{CD} . By using the same sampling argument as before, π_{CS} and V_{CS} follow the correct distribution.

It remains to show that $\widetilde{\text{EO.CT}}$ are distributed correctly. We use a hybrid argument to show this.

Hybrid₀: The first hybrid defines $\widetilde{\text{EO.CT}}$ according to the real execution. In the real execution, Party C uses the EO.Rnd procedure to rerandomize $\{\{\text{EO.ct}'_{t,i,j}\}_{j \in [m_{t,i}], i \in [m_t], t \in [T]}$ and applies the permutations π_{CD} and π_{CS} the outcome is $\widetilde{\text{EO.CT}}$.

Hybrid₁ This hybrid generates EO.CT as a fresh encryption of $p_{t,i,j}$ using $\text{EO.Enc}(\text{EO.pk}, \text{EO.pf}, p_{t,i,j})$.

Hybrid₂: The last hybrid generates EO.CT as an encryption of 0 using $\text{EO.Enc}(\text{EO.pk}, \text{EO.pf}, 0)$.

Based on the statistical ciphertext indistinguishability of EO, Hybrid₀ and Hybrid₁ generate up to negligible probability the same distribution. We can use a standard hybrid argument to show this. Let ϵ be the distinguishing probability between Hybrid₀ and Hybrid₁ and let $N = \sum_{i=1}^M m_i$ be the amount of ciphertexts, then the statistical ciphertext indistinguishability can be broken with probability $\frac{\epsilon}{N}$.

We show now that Hybrid₁ and Hybrid₂ are computationally indistinguishable based on the ciphertext indistinguishability (for secret key or evaluation key owner). The two notions give the adversary access to either EO.sk or EO.ek. Since the corrupted parties among P_1 to P_T as well as Party S do not have access to either of the keys, a weaker notion suffices in which no access to EO.sk, EO.ek is given. This weaker notion is implied by both of the ciphertext indistinguishability notions of an EO scheme.

We use a standard hybrid in which we replace step by step one of the N ciphertexts with an encryption with 0. The last hybrid matches Hybrid₂ and the first hybrid Hybrid₁. For each step we use a reduction to the ciphertext indistinguishability game in which given EO.pk, EO.pf, we need to construct a distinguisher D' that distinguishes between an encryption of $x_0 = p_{t,i,j}$ and $x_1 = 0$. We construct this distinguisher by invoking the distinguisher D between two intermediate hybrids. D' forwards EO.pk and EO.pf, it generates the view of the corrupted parties as specified by the hybrids with the exception of the one ciphertext that is different in the hybrids. D' uses the challenge ciphertext for this ciphertext. Finally D' outputs the output of D .

If D successfully distinguishes two intermediate hybrids, D' breaks the ciphertext indistinguishability for the secret key and evaluation key owner of the EO scheme. Let ϵ be an upper bound on the distinguishing probability in the ciphertext indistinguishability game. Then the distinguishing probability between Hybrid₁ and Hybrid₂ is upper bounded by ϵ/N .

The indistinguishability between Hybrid₀, Hybrid₁, and Hybrid₂ concludes our claim. \square

\square

D SECURITY ANALYSIS

D.1 Security Analysis of DPMC

PROOF. We prove Theorem 10 by proving the following two claims.

CLAIM 39. *Let the secret key encryption and the PKE scheme be IND-CPA secure, the KEM simulatable and the DDH assumption hold.*

Then there exists a simulator that generates the joint view of Party C and any subset of parties P_1 to P_T that is computationally indistinguishable from the real view.

PROOF. The joint view of Party C and the subset of corrupted parties among P_1 to P_T , identified by $\mathbb{C} \subseteq [T]$ can be generated from their inputs $KV_C, KV_{t \in \mathbb{C}}$, the outputs SH_C , and the messages $\{cta_t, ctb_t, \{\{ha_{t,i,j}\}_{j \in [m_{t,i}]},$

$\{ctc_{t,i}\}_{i \in [m_{t,i}]\}_{t \in [T]}, \{KEM.cp_{i,t}\}_{i \in [m_C], t \in [T]}\}.$

The simulator on input $KV_C, \{KV_t\}_{t \in \mathbb{C}}$, and SH_C simulates the messages as follows. It samples $(KEM.pk, KEM.sk) \leftarrow KEM.KG(1^K)$ and uses $KEM.Sim$ on input $KEM.sk, SH_C$ to compute message $\{KEM.cp_{i,t}\}_{i \in [m_C], t \in [T]}$. For all $t \notin \mathbb{C}$, it samples $sk_t \leftarrow SKE.KG(1^K)$, $cta_t \leftarrow PKE.Enc(pk_D, 0)$, $ctb_t \leftarrow SKE.Enc(sk_t, 0)$, $ctc_{t,i} \leftarrow SKE.Enc(sk_t, 0)$, $r_{t,i,j} \xleftarrow{R} \mathbb{Z}_q$ and defines $ha_{t,i,j} := g^{r_{t,i,j}}$.

We use the following sequence of hybrids to show that the joint view during the real execution is indistinguishable from the view generated by the simulator.

Hybrid₁: Identical to the view during the real protocol execution.

Hybrid₂: Computes $(KEM.pk, KEM.sk) \leftarrow KEM.KG(1^K)$ as output of $KEM.Sim$ on input $KEM.sk, SH_C$.

Hybrid₃: For all $t \in \mathbb{C}$, compute cta_t as $cta_t \leftarrow PKE.Enc(pk_D, 0)$.

Hybrid₄: For all $t \in \mathbb{C}$, compute $ctb_t, ctc_{t,i}$ as $ctb_t \leftarrow SKE.Enc(sk_t, 0)$, $ctc_{t,i} \leftarrow SKE.Enc(sk_t, 0)$.

Hybrid₅: For all $t \in \mathbb{C}$, compute $ha_{t,i,j}$ as $ha_{t,i,j} := g^{r_{t,i,j}}$ where $r_{t,i,j} \xleftarrow{R} \mathbb{Z}_q$.

Hybrid₁ and Hybrid₂ are indistinguishable except with negligible probability based on the simulatability of the key encapsulation scheme.

Hybrid₂ and Hybrid₃ are indistinguishable based on the IND-CPA security of the PKE scheme. Notice that only party D has access to sk_D . The reduction works as follows. Let there be a distinguisher against Hybrid₂ and Hybrid₃ with probability ϵ . Then, we define a sequence of $T + 1$ hybrids in which we step by step replace cta_t with encryptions of 0. The distinguisher can distinguish at least one of the hybrids with at least probability ϵ/T . We can use it to construct a distinguisher against the IND-CPA game as follows. The distinguisher receives pk from the IND-CPA game and defines $pk_D := pk$. It sets $x_0 := sk_t$ and $x_1 := 0$ and receives back a challenge ciphertext ct . It defines $cta_t := ct$. It outputs whatever the distinguisher between Hybrid₂ and Hybrid₃ outputs. This distinguisher breaks the IND-CPA security with probability ϵ/T . By the security of the PKE scheme, this must be negligible and therefore Hybrid₂ and Hybrid₃ can be distinguished with at most negligible probability as well.

Since cta is independent of the symmetric key, we can now use the IND-CPA security of the symmetric key encryption to replace ctb and ctc with encryptions of 0. Again, we define a sequence of hybrids in which we replace step by step the ciphertexts by encryptions of 0. The distinguisher against Hybrid₃ and Hybrid₄ can distinguish at least two consecutive intermediate hybrids with at least probability $\epsilon/(T + \sum_{t \in \mathbb{C}} m_t)$. The distinguisher against the IND-CPA game can use $x_0 := a_t$ ($x_0 := sh_{D,t,i}$) and x_1 in the IND-CPA game for the challenge ciphertext. Then, the distinguisher can use the challenge ciphertext to either simulate the first or second consecutive intermediate hybrid and output whatever the distinguisher against the hybrids outputs. Therefore, Hybrid₃ and Hybrid₄ can be distinguished at most with negligible probability.

Notice that the ciphertexts are now independent of scalar a_t . We can use the DDH assumption (Def. 14) to argue that Hybrid₄ and Hybrid₅ are indistinguishable. Again, we use a sequence of hybrids in which we replace step by step $ha_{t,i,j}$ with a uniform group element, i.e., $ha_{t,i,j} := g^{r_{t,i,j}}$ where $r_{t,i,j} \xleftarrow{R} \mathbb{Z}_q$. There are $\sum_{t \in \mathbb{C}, i \in [m_t]} m_{t,i}$ hybrids. Let there be a distinguisher between Hybrid₄ and Hybrid₅ with probability ϵ . Then, there are two consecutive intermediate hybrids that this distinguisher distinguishes with at least probability $\epsilon/(\sum_{t \in \mathbb{C}, i \in [m_t]} m_{t,i})$. The reduction to DDH works as follows. The DDH distinguisher receives challenge A, B, C . Before invoking the hybrid distinguisher, it programs $H_{\mathbb{G}}(p_{t,i,j}) := B$ and defines $h_{t,i,j} := C$. For all other $h_{t,i,j}$ that are not uniform yet, it programs $H_{\mathbb{G}}(p_{t,i,j}) := g^{x_{t,i,j}}$, where $x_{t,i,j} \xleftarrow{R} \mathbb{Z}_q$ and defines $h_{t,i,j} := A^{x_{t,i,j}}$. The DDH distinguisher outputs the output of the hybrid distinguisher. When $A = g^a, B = g^b, C = g^{ab}$, all $h_{t,i,j}$ are correctly defined as in the first consecutive hybrid. When A, B, C are

uniform group elements, $h_{t,i,j} = C$ is uniform while all other $h_{t,i,j}$ are distributed according to the second (and first) of the consecutive hybrids. Therefore, Hybrid₄ and Hybrid₅ can be distinguished with at most negligible probability which concludes the proof of our claim. \square

CLAIM 40. *Let the KEM scheme be key indistinguishable and the DDH assumption hold.*

Then there exists a simulator with access to the leakage defined in Def. 9 that generates the joint view of Party C and any subset of parties P_1 to P_T that is computationally indistinguishable from the real view.

PROOF. The joint view of Party D and the subset of corrupted parties among P_1 to P_T , i.e., defined by $\mathbb{C} \subseteq [T]$ can be generated by the inputs sk_D , $\{KV_t\}_{t \in \mathbb{C}}$, output SH_D and messages $KEM.pk$, $\{(h_{C,i,j})_{j \in [m_{C,i}]} \}_{i \in [m_C]}$ and $\{cta_t, ctb_t, \{(h_{C,t,i,j})_{j \in [m_{t,i}]} \}, ctc_{t,i}\}_{i \in [m_t]}\}_{t \in [T]}$.

Given the leakage defined in Def. 9 and inputs sk_D , $\{KV_t\}_{t \in \mathbb{C}}$, SH_D , the simulator works as follows. The simulator uses the leakage to define $\{(h_{C,i,j})_{j \in [m_{C,i}]} \}_{i \in [m_C]}$ and $\{hc_{t,i,j}\}_{j \in [m_{t,i}], i \in [m_t], t \in [T]}$. For all $t \notin \mathbb{C}$, it samples $a_t \xleftarrow{R} \mathbb{Z}_q$ (for all other t , a_t is already defined when generating the view for P_t). $h_{C,t,i,j} := hc_{t,i,j}^{a_t}$. For all $t \notin \mathbb{C}$, sample $sk_t \leftarrow SKE.KG(1^K)$ and use sk_t and pk_D to define cta_t, ctb_t and ctc according to the protocol description, where $sh_{D,t,i}$ is defined s.t. it is consistent with SH_D and $(KEM.cp_{t,i}, KEM.k_{t,i}) \leftarrow KEM.Enc(KEM.pk)$.

We prove that the view generated by the simulator is indistinguishable from the real view using the following hybrids.

Hybrid₁: Is identical to the view during the protocol.

Hybrid₂: For all $t \notin \mathbb{C}$, generate $(KEM.cp_{t,i}, KEM.k_{t,i}) \leftarrow KEM.Enc(KEM.pk)$. (Now $KEM.k_{t,i}$ is independent of sh_C).

Hybrid₃: Use the leakage to define $\{(h_{C,i,j})_{j \in [m_{C,i}]} \}_{i \in [m_C]}$ and $\{hc_{t,i,j}\}_{j \in [m_{t,i}], i \in [m_t], t \in [T]}$.

Hybrid₁ and Hybrid₂ are indistinguishable based on the key indistinguishability of the key encapsulation. To show this, we use a sequence of hybrids in which we replace $KEM.cp_{t,i}$ generated by P_t for $t \notin \mathbb{C}$ and related to $sh_{C,t,i}$ with $(KEM.cp'_{t,i}, KEM.k'_{t,i}) \leftarrow KEM.Enc(KEM.pk)$. We use the triangular inequality which implies that if $KEM.cp, KEM.k$ cannot be distinguished with more than probability ϵ from $KEM.cp, u$ for a uniform u , $KEM.cp, KEM.k$ cannot be distinguished from $KEM.cp', KEM.k$ with more than probability 2ϵ . Let there be a distinguisher that distinguishes Hybrid₁ and Hybrid₂ with probability ϵ . Then it distinguishes at least two consecutive intermediate hybrids with probability $\epsilon/(\sum_{t \in \mathbb{C}} m_t)$. Given this distinguisher, we build a distinguisher against the key indistinguishability which receives challenge $KEM.cp, KEM.k$ and sets $sh_{C,t,i} := KEM.k$. The distinguisher outputs the output of the hybrid distinguisher. When $KEM.k$ is consistent with $KEM.cp$, the distinguisher simulates Hybrid₁ and otherwise Hybrid₂. This distinguisher breaks the key indistinguishability with probability $\epsilon/(\sum_{t \in \mathbb{C}} m_t)$. Since this is negligible, Hybrid₁ and Hybrid₂ cannot be distinguished except negligible probability.

Hybrid₂ and Hybrid₃ are indistinguishable based on the DDH assumption. We show this by using a sequence of intermediate hybrids in which we replace $\{(h_{C,i,j})_{j \in [m_{C,i}]} \}_{i \in [m_C]}$ and $\{hc_{t,i,j}$

$\}_{j \in [m_{t,i}], i \in [m_t], t \in [T]}$ with uniform group elements. If there is a distinguisher that distinguishes Hybrid₂ and Hybrid₃ with probability ϵ , then it distinguishes at least two consecutive intermediate hybrids with probability $\epsilon/(\sum_{i \in [m_C]} m_{C,i} + \sum_{t \in [T], i \in [m_t]} m_{t,i})$. The distinguisher against DDH receives A, B, C and defines $hca_{t,i,j}^{1/a_t} := C$ ($h_{C,i,j} := C$), programs $H_{\mathbb{G}}(p_{t,i,j}) := B$ ($H_{\mathbb{G}}(c_{i,j}) := B$). For all $hca_{t,i,j}, hc_{t,i,j}$ that are not uniform yet, program $H_{\mathbb{G}}(p_{t,i,j}) := g^{x_{t,i,j}}$, $H_{\mathbb{G}}(c_{i,j}) := g^{x_{i,j}}$ and define $hca_{t,i,j}^{1/a_t} := A^{x_{t,i,j}}$, $hc_{t,i,j} := A^{x_{i,j}}$. When $A = g^a, B = g^b, C = g^{ab}$, the DDH distinguisher simulates the first of the intermediate hybrids otherwise the second one. Notice that in the latter case, $hc_{t,i,j} := hca_{t,i,j}^{1/a_t} (h_{C,i,j})$ is uniform. This concludes the proof of our claim. \square

D.2 Security Analysis of D_sPMC

PROOF. We prove Theorem 13 by constructing a simulator that can generate a view of the corrupted parties from their inputs and outputs that is indistinguishable from their view during a real execution. We emphasize that the distinguisher has access to the inputs and outputs of the honest parties specified by the ideal functionality in Fig. 2, which matches the outputs of the real protocol. We show this in the following three claims.

CLAIM 41. *Let PKE be an IND-CPA secure and correct PKE scheme, PRG a secure pseudorandom generator and EO be a correct and satisfy statistical rerandomized ciphertext indistinguishability, the (semi-honest) ciphertext indistinguishability for the evaluation key owner and ciphertext well-formedness.*

Then, there exists a simulator that generates the joint view of Party C and any subset of parties P_1 to P_T that is indistinguishable from the joint view during the protocol execution.

PROOF. The joint view can be generated from the input and messages received by party C and the subset of parties P_1 to P_T . Let this subset be $\mathbb{C} \subseteq [T]$. Notice that the parties do not have any outputs as specified in the ideal functionality in Fig. 2.

The inputs are KV_C and $\{KV_t\}_{t \in \mathbb{C}}$ and the output is SH_C . The parties P_1 to P_T receive messages $EO.pk, EO.pf$ and have access to pk_D , where $EO.pf$ is generated by Party C. Party C receives the messages $\{\{EO.ct_{t,i,j}\}_{j \in [m_{t,i}]}, sh_{C,t,i}\}_{i \in [m_t]}, cta_t\}_{t \in [T]}$ and $\{\widehat{KEM}.cp_{i,t}\}_{i \in [m_C], t \in [T]}$.

The simulator receives input $KV_C, \{KV_t\}_{t \in \mathbb{C}}, SH_C$ and emulates the view as follows. It samples $(KEM.pk, KEM.sk) \leftarrow KEM.KG(1^K)$, $(EO.pk, EO.sk) \leftarrow EO.KG(1^K)$ and $(pk_D, sk_D) \leftarrow PKE.KG(1^K)$. It samples $cta_t \leftarrow PKE.Enc(pk, 0)$ for all $t \notin [T]$. It samples $sh_{C,t,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ for all $t \notin \mathbb{C}$. It defines $\widehat{sh}_{C,i,t}$ consistently with SH_C for all $t \in [T]$ and defines $\widehat{KEM}.cp_{i,t} \leftarrow KEM.Sim(KEM.sk, \widehat{sh}_{C,i,t})$. Further, it defines $EO.ct_{t,i,j} \leftarrow EO.Enc(EO.pk, EO.pf, 0)$ for all $t \notin \mathbb{C}$. For $t \in \mathbb{C}$, $EO.ct_{t,i,j} \leftarrow EO.Enc(EO.pk, EO.pf, p_{t,i,j})$, where $p_{t,i,j} \in KV_t$.

We use the following sequence of hybrids to show that the simulated view is indistinguishable from the view during the real protocol execution.

Hybrid₀: Identical to the view during the real protocol execution.

Hybrid₁: Samples $cta_t \xleftarrow{PKE.Enc} (pk, 0)$ for all $t \notin \mathbb{C}$.

Hybrid₂: Samples $sh_{D,t,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ for all $t \notin \mathbb{C}$ (instead of using PRG).

Hybrid₃: Invoke the simulator of the shuffling protocol to simulate the view during the shuffling. The input $\{\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}, sh_{C,t,i}\}_{i \in [m_t]}, \{\widetilde{\text{EO.ct}}_{i,j}\}_{i \in [M], j \in [m_i]}$ of the simulator is distributed as in Hybrid₂. Notice that the simulator also receives EO.pk, EO.pf and EO.ek.

Hybrid₄: Replaces $\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}$ for all $t \notin \mathbb{C}$ and all $\{\widetilde{\text{EO.ct}}_{i,j}\}_{i \in [M], j \in [m_i]}$ with independent encryptions of 0. More precisely, $\text{EO.ct}_{t,i,j} \leftarrow \text{EO.Enc}(\text{EO.pk}, \text{EO.pf}, 0)$ and $\widetilde{\text{EO.ct}}_{i,j} \leftarrow \text{EO.Enc}(\text{EO.pk}, \text{EO.pf}, 0)$.

Hybrid₅: Samples $sh_{C,t,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ for all $t \notin \mathbb{C}$. Further, defines $\widehat{sh}_{C,i,t}$ consistently with SH_C and samples $\widehat{\text{KEM.cp}}_{i,t} \leftarrow \text{KEM.Sim}(\text{KEM.sk}, \widehat{sh}_{C,i,t})$.

Notice that the view in Hybrid₅ is identically distributed as the view generated by the simulator.

We now show that the hybrids are indistinguishable. Let Hybrid₀ and Hybrid₁ be distinguishable with probability ϵ . We define a sequence of intermediate hybrids that replaces the ciphertexts $\text{cta}_t \leftarrow \text{PKE.Enc}(\text{pk}, \text{seed}_t)$ with $\text{cta}_t \xleftarrow{\text{PKE.Enc}} (\text{pk}, 0)$. Then there is a distinguisher that distinguishes one of the intermediate hybrids with at least probability ϵ/τ . Such a distinguisher would directly distinguish challenge ciphertexts for $x_0 := \text{seed}_t$ from $x_1 := 0$ in the IND-CPA game of the PKE scheme. Therefore the distinguishing probability between Hybrid₀ and Hybrid₁ is upper bounded by the IND-CPA security of PKE.

Let Hybrid₁ and Hybrid₂ be distinguishable with probability ϵ . We define a sequence of intermediate hybrids in which we step by step replace $(sh_{D,t,1}, \dots, sh_{D,t,m_t}) = \text{PRG}(\text{seed}_t)$ with $(sh_{D,t,1}, \dots, sh_{D,t,m_t}) \leftarrow \{0, 1\}^{m_t \cdot |\mathcal{V}|}$. Then, there would be a distinguisher that distinguishes two consecutive intermediate hybrids with at least probability ϵ/τ . This would imply a distinguisher that breaks the security of the PRG with the same probability. Since the PRG is indistinguishable except negligible probability, Hybrid₀ and Hybrid₁ cannot be distinguished except negligible probability.

Let Hybrid₂ and Hybrid₃ be distinguishable with probability ϵ . Then, this would allow to distinguish the simulated view during the shuffle protocol from the real view. However, as shown in Theorem 35 this probability is upper bounded the correctness, the statistical rerandomized ciphertext indistinguishability, the (semi-honest) ciphertext indistinguishability (for evaluation key or secret key owner) and ciphertext well-formedness of the EO scheme. Therefore Hybrid₂ and Hybrid₃ cannot be distinguished beyond the bound given in the proof of Theorem 35.

We use the (semi-honest) ciphertext indistinguishability for the evaluation key owner to argue that Hybrid₃ and Hybrid₄ are indistinguishable. Notice that in the ideal shuffle functionality (see Fig. 12), the ciphertext sets $\{\text{EO.ct}_{t,i,j}\}_{j \in [m_{t,i}]}$ and $\{\widetilde{\text{EO.ct}}_{i,j}\}_{i \in [M], j \in [m_i]}$ are independent encryptions. Therefore, we can replace them independently with encryptions of 0. We need to use the ciphertext indistinguishability for the evaluation key owner since the simulator need access to EO.ek which is also used by the simulator of the shuffling protocol. The indistinguishability between Hybrid₃ and Hybrid₄ follows from a straightforward reduction to the ciphertext indistinguishability using a hybrid argument in which we

replace step by step each ciphertext with an encryption of 0 until all ciphertexts are encryptions of 0. If there exists a distinguisher between Hybrid₃ and Hybrid₄ that distinguishes them with probability ϵ , then there is a distinguisher that distinguishes one of the intermediate hybrids with at least probability $\epsilon/2N$, where N is the size of $\{\text{KV}_t\}_{t \in [T]}$. The distinguisher for the intermediate hybrids would then lead to a distinguisher against the ciphertext indistinguishability for the evaluation key owner of the EO scheme.

We finalize the claim by showing the indistinguishability of Hybrid₄ and Hybrid₅. Similar as in case of the ciphertexts, the ideal shuffle functionality samples the shares $sh_{C,t,i}$ and $\widehat{sh}_{C,i,t}$ independently. Therefore, we can also sample them independently. Hybrid₅ generates statistically the same view as Hybrid₄ for the following reason. Sampling $sh_C \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ and $sh_D \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ under the constraint that $sh_C \oplus sh_D = v$ (Hybrid₄) results in the same distribution as when sampling $sh_C \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ and defining $sh_D := v \oplus sh_C$ (Hybrid₅), where sh_D and v are not known to the simulator. Thus, sh_C can be sampled independently of sh_D and v by sampling $sh_C \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$. Further, by the property of KEM.Sim , $\widehat{\text{KEM.cp}}_{i,t}$ has the same distribution when being an output of KEM.Enc and KEM.Sim . This concludes our claim. \square

CLAIM 42. *Let PKE be a correct PKE scheme, KEM a secure and correct key encapsulation scheme and EO secure, correct and evaluated ciphertext simulatable.*

Then, there exists a simulator with access to the leakage graph of Def. 12 that generates the joint view of Party D and any subset of parties P_1 to P_T that is indistinguishable from the joint view during the protocol execution.

PROOF. The joint view of Party D and the subset of parties P_1 to P_T (defined by $\mathbb{C} \subseteq [T]$) can be generated from inputs $(\text{pk}_D, \text{sk}_D)$, $\{\text{KV}_t\}_{t \in \mathbb{C}}$, output SH_D and messages KEM.pk , EO.pf , $\{\text{cta}_t\}_{t \in \mathbb{C}}$ and $\{\text{KEM.cp}_i, \widehat{sh}_{D,i}\}_{i \in [M]}, \{h_{C,i,j}\}_{i \in [m_C], j \in [m_{C,i}]}, \{\text{EO.ect}_{i,j}\}_{i \in [M], j \in [m_i]}$. Further, the view depends on the leakage graph defined in Def. 12.

The simulator emulates the joint views as follows. It samples $(\text{KEM.pk}, \text{KEM.sk}) \leftarrow \text{KEM.KG}(1^\kappa)$ and $(\text{EO.pf}, \text{EO.ek}) \leftarrow \text{EO.EKG}$. The simulator defines $\{h_{C,i,j}\}_{i \in [m_C], j \in [m_{C,i}]}$ and $\{h_{i,j}\}_{i \in [M], j \in [m_i]}$ such that they are consistent with the leakage graph. It then defines $\text{EO.ect}_{i,j} \leftarrow \text{EO.Sim}(\text{EO.pk}, \text{EO.sk}, h_{i,j})$. It samples $\widetilde{sh}_{D,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ and $(\text{KEM.cp}_i, \text{KEM.k}_i) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$. Define $\widetilde{sh}_{D,i}$ s.t. that it is consistent with SH_D and $\widetilde{sh}_{D,i}$. For all $\widetilde{sh}_{D,i}$ that are not defined yet, sample $\widetilde{sh}_{D,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$. For $t \notin \mathbb{C}$, it samples $\text{seed}_t \xleftarrow{R} \{0, 1\}^\kappa$ and $\text{cta}_t \leftarrow \text{PKE.Enc}(\text{pk}, \text{seed}_t)$, which is identical to the protocol description.

We prove the claim by using the following sequence of hybrids.

Hybrid₀: Is identical to the views during the real execution of the protocol.

Hybrid₁: Simulates the view during the shuffling by using the simulator of the shuffle protocol.

Hybrid₂: Sample $\text{EO.ect}_{i,j} \leftarrow \text{EO.Sim}(\text{EO.pk}, \text{EO.sk}, h_{i,j})$, where $h_{i,j} := \text{EO.Eval}(\text{EO.ek}, p_{i,j})$ and $p_{i,j}$ is the reshuffled $p_{t,i,j}$, which can be computed from the shuffle permutation π and $\{\text{KV}_t\}_{t \in [T]}$.

Hybrid₃: It defines $\{h_{C,i,j}\}_{i \in [m_C], j \in [m_{C,i}]}$ and $\{h_{i,j}\}_{i \in [M], j \in [m_i]}$ such that they are consistent with the leakage graph.

Hybrid₄: Samples $(\text{KEM.cp}_i, \text{KEM.k}_i) \leftarrow \text{KEM.Enc}(\text{KEM.pk})$ s.t. it is independent of $\text{sh}_{D,i}$, $\tilde{\text{sh}}_{D,i}$ and SH_D .

Hybrid₅: Samples $\tilde{\text{sh}}_{D,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$ and defines $\tilde{\text{sh}}_{D,i}$ s.t. that it is consistent with $\text{SH}_{\mathcal{J},D}$ and $\tilde{\text{sh}}_{D,i}$. For all $\tilde{\text{sh}}_{D,i}$ that are not defined yet, sample $\tilde{\text{sh}}_{D,i} \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$.

If Hybrid₀ and Hybrid₁ can be distinguished with probability ϵ , then there is a distinguisher against the simulator of the shuffle protocol with probability ϵ . Since such a distinguishing probability is negligible (based on the security of EO, see Theorem 35), distinguishing Hybrid₀ from Hybrid₁ is also negligible.

If Hybrid₁ and Hybrid₂ can be distinguished with probability ϵ , we can define a sequence of hybrids that step by step replaces $\text{EO.ect}_{i,j}$ with outputs of EO.Sim . Now, there are at least two consecutive intermediate hybrids that can be distinguished with probability $\epsilon/(\sum_{i=1}^M m_i)$. Since this probability is negligible due to the evaluated ciphertext simulatability of EO, Hybrid₁ and Hybrid₂ can also only be distinguished with negligible probability.

In Hybrid₂ $h_{i,j}$ and $h_{C,i,j}$ are the outputs of EO.Eval whereas in Hybrid₃ they are uniform in $\{0, 1\}^K$. We prove that Hybrid₂ and Hybrid₃ are indistinguishable except with negligible probability by a reduction to the pseudorandomness of EO.Eval . Let there be a distinguisher distinguishing Hybrid₂ and Hybrid₃ with probability ϵ , then we can build a distinguisher against the pseudorandomness of EO.Eval with probability ϵ . The latter requests all $h_{i,j}$ and $h_{C,i,j}$ from the EO.Eval oracle, uses them to simulate Hybrid₂, Hybrid₃ and outputs the output of the former distinguisher. When they are actual EO.Eval outputs, it simulates Hybrid₂ and when they are uniform, it simulates Hybrid₃.

If Hybrid₃ and Hybrid₄ can be distinguished with probability ϵ , we can define a sequence of hybrids that step by step replaces $(\text{KEM.cp}_i, \text{KEM.k}_i)$ with $(\text{KEM.cp}_i, \text{KEM.k}'_i)$ where $(\text{KEM.cp}_i, \text{KEM.k}_i) \leftarrow \text{KEM.Enc}(\text{pk})$. Now there exist two consecutive intermediate hybrids that can be distinguished which implies a distinguisher for $(\text{KEM.cp}, \text{KEM.k})$ and $\text{KEM.cp}, \text{KEM.k}'$ with probability ϵ/m . By the triangular inequality, we can then build a distinguisher for $(\text{KEM.cp}, \text{KEM.k})$ and $(\text{KEM.cp}, u)$ with probability $\epsilon/2$, where $u \xleftarrow{R} \{0, 1\}^{|\mathcal{V}|}$. Such a distinguisher breaks the key indistinguishability for the KEM. Since this is negligible, Hybrid₃ and Hybrid₄ cannot be distinguished except with negligible probability.

Hybrid₄ and Hybrid₅ produce identically distributed views. Note that $\tilde{\text{sh}}_{D,i}$, $\text{sh}_{D,i}$ and $\text{SH}_{\mathcal{J},D}$ are independent of $(\text{KEM.cp}_i, \text{KEM.k}_i)$.

Further, due to the simulator of the shuffling, they are independent of $\text{sh}_{C,t,i}$ and $\text{sh}_{D,t,i}$. Therefore, they can be sampled independently which concludes the proof of our claim. \square

CLAIM 43. *Let EO be a secure and correct randomizable encrypted OPRF scheme. Then, there exists a simulator that generates the joint view of Party S and any subset of parties P_1 to P_T that is indistinguishable from the joint view during the protocol execution.*

PROOF. The joint view of Party S and the corrupted subset of parties P_1 to P_T (defined by set $\mathbb{C} \subseteq [T]$) can be generated from their input $\{\text{KV}_t\}_{t \in \mathbb{C}}$ and the received messages $\{\text{sh}_{C,i}\}_{i \in [M]}$ and KEM.pk .

The simulator samples $\tilde{\text{sh}}_{C,i} \xleftarrow{R} \{0, 1\}$ and uses the simulator of the shuffle protocol to simulate the view during the shuffling.

The view generated by the simulator is indistinguishable from the view during the real protocols by the indistinguishability of the simulated view of shuffling from the real view of the shuffling. Notice that in case of using the simulated view of the shuffling, $\{\tilde{\text{sh}}_{C,i}\}_{i \in [M]}$ are independent of $\{\text{sh}_{C,t,i}\}_{i \in [m_t], t \in [T]}$. Therefore, $\tilde{\text{sh}}_{C,i}$ can be sampled independently when using the simulated view during the shuffling. \square

\square

E EXTENDING LEFT JOIN TO INNER JOIN

DPMC and D_s PMC can be extended to support other types of joins such as an inner join instead of a left join. In both protocols, party D performs the join based on the encrypted datasets of C and all delegators (i.e., in DPMC in step ④ and in D_s PMC in step ⑨). Performing the left join in party D hides from party C which of its rows have been matched with one of the delegators' rows and which have not. It is straightforward to extend our delegated protocols to compute the inner join (i.e., intersection) between KV_C and KV_P and secret share the associated metadata for these rows. This can be performed very efficiently using hash join over the encrypted identifiers and sending the KEM.cp value to C only for the records present in both datasets. Notably, computing the inner join leaks the intersection size to party C but also renders the downstream MPC computation more efficient since it does not have to process secret shares of NULL.