

# A Proposal for Assisted Private Information Retrieval

Natnatee Dokmai<sup>1</sup>, L. Jean Camp<sup>1</sup>, and Ryan Henry<sup>2</sup>

<sup>1</sup>Indiana University, Bloomington

<sup>2</sup>University of Calgary

## Abstract

Private Information Retrieval (PIR) addresses the cryptographic problem of hiding sensitive database queries from database operators. In practice, PIR schemes suffer from either high computational costs or from restrictive requirements difficult to justify in practical settings. In this work, we introduce *Assisted Private Information Retrieval* (APIR), a new PIR problem for keyword-value databases which generalizes and relaxes the database consistency assumption in multi-server PIR. Leveraging the decentralized nature of Domain Name Service (DNS), APIR is able to address a privacy issue inherent to encrypted DNS proposals such as DNS-over-HTTPS (DoH) by preventing DNS operators from collecting sensitive data. We propose a construction of *Synchronized APIR*, an efficient hybrid APIR scheme between black-box single-server PIR and non-black-box multi-server PIR. We apply Synchronized APIR to a proof-of-concept protocol for private DNS query, and demonstrate that APIR is able to outperform the baseline single-server PIR protocol after the initial one-time cost.

## 1 Introduction

One often overlooked privacy aspect of information access over a network is the monitoring of information queries by database-operating servers. To illuminate this concern, let us consider a scenario: a user wishes to download a sensitive file from a website, but the act of access to that particular file is incriminating. At this point, the user has a few options. She can use an encrypted channel like TLS to prevent eavesdropping, but this does not prevent the leakage of the query to the web server; or she can use anonymizing technology like VPN or Tor to hide her identity, but this does not prevent the web server from being notified and collecting statistics about file access.

A case in point is the Domain Name System (DNS). DNS privacy and security are the core argument for encrypted DNS such as DNS-over-HTTPS (DoH) and DNS-over-TLS (DoT). However, the present encrypted DNS

proposals do not necessarily increase consumers' privacy. The most notable privacy risk is increased data concentration at DNS operators. Given the exploitation of users' data by data collectors and online services in the recent years, it is not far-fetched to claim that concentration of data is a privacy violation.

The family of cryptographic protocols that proposes to address this issue is *private information retrieval* (PIR). The core privacy definition of PIR requires that any two different PIR queries appear indistinguishable from the point of view of the attacker, thus preventing database operators from collecting sensitive information about the queries. There have been many proposed PIR schemes in recent years, yet they all suffer from similar practical issues. The PIR schemes that require one server to operate, or *single-server* PIR, are applicable to practical scenarios; yet, they often rely on additively homomorphic encryption schemes which are computationally expensive and incur high communication costs in practice. The PIR schemes that require multiple servers to operate, or *multi-server* PIR, are orders of magnitude cheaper than single-server PIR to deploy, yet they have conflicting requirements that 1) all the servers hold some form of the identical copy of the database, and 2) some servers do not collude. In practice, it is difficult to imagine a scenario in which both requirements would hold simultaneously. In the DNS application for instance, it can be argued that DNS servers administered by independent organizations are unlikely to collude. However, to require that they hold identical copies cache tables and zone files seem impossible. Further, it makes the application less useful e.g. it does not allow the same domain to point to different local IP addresses, as in the case of scalable cloud services.

### 1.1 Our contributions

In this work, we introduce *assisted private information retrieval* (APIR). APIR generalizes multi-server PIR in the following ways, while still maintaining the non-collusion assumption:

1. The databases are keyword-value maps instead of indexed vectors.

- Client can query the *main server* with assistance from *assisting servers* without requiring the assisting servers to hold an identical copy of the database.

We present *Synchronized APIR*, an APIR scheme which allows the client to synchronize the “view” of the databases across the servers prior to making queries. Synchronized APIR is a hybrid protocol between a black-box single-server PIR scheme and a non-black-box multi-server PIR scheme similar to CGKS [7]. The hybridization takes advantage of the cheaper multi-server PIR scheme to “assist” in lowering the costs of the single-server PIR scheme between the client and main server; the level of assistance is determined by the amount of overlap between main server’s and assisting servers’ databases. Because Synchronized APIR makes use of a black-box single-server scheme, improvements in single-server PIR in the future will also lead to improvements in Synchronized APIR. We provide a formal analysis and an implementation of Synchronized APIR in Rust using SealPIR [4] as the underlying single-server scheme.

We apply Synchronized APIR to demonstrate a proof-of-concept private DNS query application, specifically to query nameserver (NS) records, among DNS cache servers. Then, we evaluate the application with simulated datasets based on realistic assumptions about DNS queries and cache behavior.

The results show that despite the higher initial one-time cost, private DNS query via Synchronized APIR is able to outperform SealPIR, the baseline single-server PIR, either in communication cost by a factor of 5 or in computational cost by a factor of 8, given the most optimal popularity distribution of DNS queries.

## 2 Related Works

Our work is partly inspired by Fanti *et al.* [9], who propose an information-theoretic multi-server PIR scheme over unsynchronized databases. In this setting, the PIR servers hold different copies of the same indexed database *with some database values missing*. The scheme proceeds by first allowing the client to synchronize the database “view” to identify the missing values. Then, the client makes a query in a way that avoids the missing values, while also hide which values are missing from the servers. This setting is applicable to peer-to-peer (P2P) file-sharing, where P2P users share and download parts of the same file or database in small, indexed fragments.

We view Fanti *et al.*’s scheme as a solution to a subset of the APIR problem with two additional limitations: 1) the database is strictly indexed, and 2) the client is not allowed to query any values with missingness. (Indeed, Synchronized APIR also directly solves this problem.) However, we remark that the second limitation

may be counter-productive to privacy in practice. Sensitive data is often the rarer one and the one that goes missing; this introduces a contradiction to Fanti *et al.*’s scheme where the more sensitive the data is, the more likely it is unavailable. Clients therefore must resort to a less privacy-preserving service to be able to access the data. APIR insists that all database values are retrievable for this reason, and Synchronized APIR solves this problem by making a separate single-server PIR query to the values with missingness.

## 3 Preliminaries

### 3.1 Single-Server PIR

Single-server PIR, or sPIR for short, is a PIR protocol that requires only one database-operating server to interact with the client to query information. sPIR is associated with *computational* PIR (CPIR) because of the computational assumptions usually required to instantiate this scheme. One example of an sPIR scheme is the *trivial* PIR, where the client downloads the entire database from the server and picks out the value she wants to query locally. This technically satisfies the privacy requirements of PIR because the server is unable to learn which value the client wishes to query. However, this results in a high communication, which is prohibitive in practice. Therefore, any non-trivial PIR schemes must aim to satisfy the privacy requirements as well as keep the communication cost lower than the trivial PIR.

We formally define sPIR below in Definition 3.1. A typical flow of an sPIR protocol proceeds as follows. 1) The client generates the public-secret key pair with *SGEN* and gives the public key to the server. 2) Once the client has decided to query the  $i$ -th item in the database, an sPIR query is generated with *SQUERY* using  $i$ . The query is sent to the server. 3) The server generates a reply from the query and database using *SREPLY*. The reply is returned to the client. 4) The client decodes the reply with the secret key using *SDECODE* to obtain the answer.

**Definition 3.1** (sPIR Scheme). Given a Server operating an indexed database  $V = (v_0, \dots, v_{m-1})$  for all  $i \in \{0, \dots, m-1\}$ . An sPIR scheme is a tuple  $\Pi^{\text{sPIR}} = (\text{SGEN}, \text{SQUERY}, \text{SREPLY}, \text{SDECODE})$  defined by

- $\text{SGEN}(1^\lambda) \rightarrow (\text{spk}, \text{ssk})$  where  $\lambda$  is the security parameter,  $\text{spk}$  the public key, and  $\text{ssk}$  the secret key.
- $\text{SQUERY}(\text{spk}, i, m) \rightarrow \text{sq}$  generates query  $\text{sq}$  from index  $i \in \{0, \dots, m-1\}$  targeting item  $v_i$ .
- $\text{SREPLY}(\text{spk}, V, \text{sq}) \rightarrow \text{sr}$  generates reply  $\text{sr}$  for query  $\text{sq}$  from database  $V$ .

- $\text{SDECODE}(\text{ssk}, \text{sr}) \rightarrow \text{sa}$  decodes reply  $\text{sr}$  for answer  $\text{sa}$ .

The formal definition of sPIR correctness is provided in Definition A.1. Intuitively, an sPIR scheme is correct if the answer is the same as the target value in the query.

The formal definition of sPIR privacy is provided in Definition A.2. Intuitively, an sPIR scheme is privacy-preserving if the adversary is unable to distinguish between a query encoding index  $i$  and a query encoding index  $j$ , even if the adversary can choose  $i$  and  $j$  and has oracle access to SQUERY.

Indeed, there are many existing schemes in literature that can satisfy these definitions of sPIR, notably because an additively homomorphic encryption scheme implies an sPIR scheme [14]. An additively homomorphic encryption scheme can be used as a building block as follows: 1) For a given index  $i$ , the query is encoded as the  $i$ -th standard basis vector  $e_i$ ; 2) the client encrypts  $e_i$  component-wise using the homomorphic encryption scheme; 3) the server homomorphically evaluates the inner product, treating the database values as scalar constants so that the result is a linear combination; 4) the client decrypts the result. Note that this is merely a simplified construction. There have been multiple techniques introduced to lower computational and communication cost in practice, including recursive PIR [3], which treats the database as two-dimensional and accesses it by row and column index; and ciphertext packing [4], which packs multiple query indices into a single ciphertext.

Although this instantiation of sPIR seems simple enough, in practice the computational cost can be too high—sometimes to the point where it can be faster for the client to download the entire database [13]. In recent years, the development has been moving toward lattice-based cryptography [2, 4, 11], whose efficiency has enabled the PIR schemes to become more practical.

In this work, we use SealPIR [4], one of the state-of-the-art sPIR schemes, as a building block for our scheme. In addition to the relatively low computational cost utilizing the aforementioned techniques, SealPIR offers a significant advantage over its predecessors via its query compression technique, which cuts down the size of a query by a large factor.

### 3.2 Multi-Server PIR

Multi-server PIR, or mPIR for short, is a PIR scheme that requires multiple database-operating servers to interact with the client to query information. mPIR is a broader category of *information-theoretic* PIR (IT-PIR) in literature, in that all IT-PIR schemes are multi-server, but some mPIR schemes can be a hybrid between IT-PIR and CPIR. Existing mPIR schemes rely on the following two assumptions: 1) databases must be *consistent*,

meaning that each server holds the same “ground-truth” copy of the database which may be preprocessed for the scheme; and 2) no more than a given number of servers can *collude*, meaning that they cannot share the client’s private information other than what is instructed. One of the well-known mPIR schemes is the CGKS scheme [7]. Below, we provide a 3-server example of CGKS, although CGKS can naturally be extended to an unbounded number of servers.

Suppose there are three database-operating servers,  $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3$ , each holding a consistent database  $V = (a, b, c, d, f) \in \text{GF}(2)^{\ell \cdot 5}$ ; and suppose the client  $\mathcal{C}$  wishes to retrieve the item at index 2 i.e.  $c$ , then she must generate 3 queries, one for each server, following the steps below.

1. Sample  $q_1 \in \text{GF}(2)^5$  and  $q_2 \in \text{GF}(2)^5$  uniformly at random. Assume from now on that  $q_1 = (0, 0, 1, 1, 0)$  and  $q_2 = (0, 1, 0, 0, 0)$ .
2. Compute  $q_3 \leftarrow q_1 \oplus q_2$ . Therefore  $q_3 = (0, 1, 1, 1, 0)$ .
3. Encode index 2 as  $e_3 = (0, 0, 1, 0, 0)$  (or index  $i$  as  $e_{i+1}$  in general).
4. XOR  $q_3$  with  $e_3 = (0, 0, 1, 0, 0)$ , namely  $q_3 \leftarrow q_3 \oplus e_3$ . Therefore  $q_3 = (0, 1, 0, 1, 0)$ .

Next,  $\mathcal{C}$  sends each query  $q_i$  to server  $\mathcal{S}_i$ , who then generates a reply  $r_i$  by computing a dot product  $r_i \leftarrow q_i \cdot V$ . As a result, we have  $r_1 = c \oplus d$ ,  $r_2 = b$ , and  $r_3 = b \oplus d$ .  $\mathcal{S}_i$  returns reply  $r_i$  to  $\mathcal{C}$ . Finally,  $\mathcal{C}$  decodes the replies to obtain the answer by computing  $r_1 \oplus r_2 \oplus r_3 = c$ .

The CGKS scheme generalizes to any number of servers with databases of any size. However, it is worth remarking that CGKS scheme as stated above is only non-trivial (i.e. communication cost lower than downloading the entire database) if the database is not exponentially lopsided (i.e., if the length of each record is super-logarithmic in the number of records). Otherwise, the query size is asymptotically equivalent to database size.

Correctness and privacy is defined in the same manner as sPIR. The privacy definition includes the notion of collusion threshold in addition: that a group of  $t$  or fewer colluding servers cannot distinguish between two queries. This property is called *t-collusion resistance*. CGKS scheme is  $(n - 1)$ -collusion-resistant where  $n$  is the number of servers. This is because any  $n - 1$  queries are statistically indistinguishable for a uniformly random string of the same length.

In our construction that follows, we will use as a building block a variation of CGKS that does not require consistent databases and uses a PRG to generate queries in order to reduce the communication cost.

## 4 Assisted PIR

The goal of this section is to provide a generic definition of Assisted PIR (APIR) including that of privacy, correctness, and communication efficiency. We will show how to construct our APIR scheme called ‘‘Synchronized APIR’’ in the upcoming section.

APIR comprises three groups of participating parties with distinct roles:

### Notations

- $\mathcal{C}$ , client
- $\mathcal{S}_0$ , main server
- $\mathcal{S}_1 - \mathcal{S}_n$ ,  $n$  assisting servers

$\mathcal{C}$  wishes to retrieve the value associated with key  $k^*$  from keyword-value database  $\text{DB}_0$ , operated by  $\mathcal{S}_0$ , with assistance from  $\mathcal{S}_1, \dots, \mathcal{S}_n$  who independently operate keyword-value databases  $\text{DB}_1, \dots, \text{DB}_n$ , respectively.  $\text{DB}_1, \dots, \text{DB}_n$  do not need to be the exact duplicates of  $\text{DB}_0$ , although they may have some common keyword-value pairs. Note that this is in contrast to sPIR and mPIR in Section 3 where databases are assumed to be index-based and consistent. APIR is formally defined below.

### Notations

- $\text{ID} := \{0, \dots, n\}$ , the set of server ID’s
- $\text{AID} := \{1, \dots, n\}$ , the set of assisting server ID’s
- $(a_i)_I := (a_{i_1}, \dots, a_{i_n})$ ;  $i_j \in I$ , a sequence ordered by the index set  $I$ . We sometimes use  $(a_i)_{i \in I}$  instead for clarity.
- $\text{Keys}(\text{DB})$  the set of all keywords in database  $\text{DB}$ .

**Definition 4.1** (APIR Scheme). Given a set of  $n + 1$  database-operating servers and define  $\text{ID} := \{0, \dots, n\}$  the set of server ID’s. An APIR scheme is a tuple  $\Pi^{\text{APIR}} = (\text{SERVGEN}, \text{CLIGEN}, \text{QUERY}, \text{REPLY}, \text{DECODE})$  defined by

- $\text{SERVGEN}(\text{id}, \text{DB}_{\text{id}}) \rightarrow \text{par}_{\text{id}}$ , where  $\text{id} \in \text{ID}$  and  $\text{par}_{\text{id}}$  is the database parameter for  $\text{DB}_{\text{id}}$ .
- $\text{CLIGEN}(1^\lambda, t, (\text{par}_{\text{id}})_{\text{ID}}) \rightarrow (\text{pk}, \text{sk})$ , where  $\lambda$  is the security parameter,  $t$  the collusion threshold where  $1 \leq t \leq n$ ,  $\text{pk}$  is the public key, and  $\text{sk}$  the secret key.
- $\text{QUERY}(\text{pk}, \text{sk}, k) \rightarrow (\text{q}_{\text{id}})_{\text{ID}}$  generates query  $\text{q}_{\text{id}}$  for Server  $\text{id}$  from query keyword  $k \in \text{Keys}(\text{DB}_0)$  targeting database value  $\text{DB}_0[k]$ .
- $\text{REPLY}(\text{id}, \text{pk}, \text{DB}_{\text{id}}, \text{q}_{\text{id}}) \rightarrow \text{r}_{\text{id}}$  generates reply  $\text{r}_{\text{id}}$  for query  $\text{q}_{\text{id}}$  from database  $\text{DB}_{\text{id}}$ , where  $\text{id} \in \text{ID}$ .

- $\text{DECODE}(\text{sk}, (\text{r}_{\text{id}})_{\text{ID}}) \rightarrow \text{a}$  decodes replies  $(\text{r}_{\text{id}})_{\text{ID}}$  for answer  $\text{a}$ .

The flow of the scheme is similar to what we have previously seen in sPIR and mPIR in Section 3: 1) The servers generate database parameters using  $\text{SERVGEN}$  and send them to client. 2) Client generates a public-secret key pair from the database and security parameters using  $\text{CLIGEN}$ . 3) Client chooses a query keyword and generates queries using  $\text{QUERY}$ ; each server gets their own query. 4) The servers responds to the query with their own database using  $\text{REPLY}$ . And finally, 5) the client aggregates all the replies and decodes to obtain the answer with  $\text{DECODE}$  using the secret key.

The correctness of APIR is defined according to what is contained in  $\text{DB}_0$ . If the client query with keyword  $k$ , then  $\text{DB}_0[k]$  is defined to be the correct answer even though there may be other  $\text{DB}_{\text{id}}$  such that  $\text{DB}_{\text{id}}[k] \neq \text{DB}_0[k]$ ; formally

**Definition 4.2** (APIR Correctness). Following Definition 4.1, scheme  $\Pi^{\text{APIR}}$  is APIR-correct for any set  $\text{ID}$ , databases  $\text{DB}_{\text{id}}$  for all  $\text{id} \in \text{ID}$ , and  $k \in \text{Keys}(\text{DB}_0)$  if

- $\forall \text{id} \in \text{ID} : \text{par}_{\text{id}} \leftarrow \text{SERVGEN}(\text{id}, \text{DB}_{\text{id}})$
- $(\text{pk}, \text{sk}) \leftarrow \text{CLIGEN}(\lambda, t, (\text{par}_{\text{id}})_{\text{ID}})$
- $(\text{q}_{\text{id}})_{\text{ID}} \leftarrow \text{QUERY}(\text{pk}, \text{sk}, k)$
- $\forall \text{id} \in \text{ID} : \text{r}_{\text{id}} \leftarrow \text{REPLY}(\text{id}, \text{pk}, \text{DB}_{\text{id}}, \text{q}_{\text{id}})$
- $\text{a} \leftarrow \text{DECODE}(\text{sk}, (\text{r}_{\text{id}})_{\text{ID}})$

then  $\text{a} = \text{DB}_0[k]$ .

The privacy definition of APIR is similar to that of mPIR in that we want to capture both the indistinguishability of queries when collusion does not exceed a certain threshold. Intuitively, an APIR scheme is privacy preserving if the attacker, who compromises a set of servers, cannot distinguish between two queries generated from query keywords  $k_0, k_1$  of the attacker’s own choice. Moreover, we want to capture a real-world scenario in which the client is not informed of which servers are compromised, or if they are compromised at all (and thus the role of the oracle  $\mathcal{O}$  in the definition below to relay the compromised queries to the attacker without the client’s knowledge). We formally define this using a game-based definition below.

**Definition 4.3** ( $(\lambda, t)$ -APIR Privacy). Define a APIR privacy experiment  $\text{Priv}_{\mathcal{A}, \Pi^{\text{APIR}}, t}(1^\lambda)$  for an APIR scheme  $\Pi^{\text{APIR}}$  according to Definition 4.1 and adversary  $\mathcal{A}$  below.

1.  $\mathcal{A}$  chooses well-formed database parameters  $(\text{par}_{\text{id}})_{\text{ID}}$  and outputs  $(\text{par}_{\text{id}})_{\text{ID}}$ .  $\mathcal{A}$  chooses a collusion set  $C \subset \text{ID}$  such that  $|C| \leq t$  and sends  $C$  to the oracle  $\mathcal{O}$ .

2. The parameters are generated by  $(\text{pk}, \text{sk}) \leftarrow \text{CLIGEN}(\lambda, t, (\text{par}_{\text{id}})_{\text{ID}})$  and  $\text{pk}$  is given to  $\mathcal{A}$ .
3.  $\mathcal{A}$  is given oracle access to QUERY in the following way:  $\mathcal{A}$  chooses and outputs  $k \in \mathcal{K}(\text{par}_0)$  where  $\mathcal{K}(\text{par}_0)$  is the query keyword space determined by a well-formed  $\text{par}_0$ . Queries are generated by  $(\mathbf{q}_{\text{id}})_{\text{ID}} \leftarrow \text{QUERY}(\text{pk}, \text{sk}, k)$  and  $(\mathbf{q}_{\text{id}})_{\text{ID}}$  is given to  $\mathcal{O}$ .  $\mathcal{O}$  gives  $(\mathbf{q}_{\text{id}})_C$  to  $\mathcal{A}$ .
4.  $\mathcal{A}$  chooses  $k_0^*, k_1^* \in \mathcal{K}(\text{par}_0)$ , and outputs  $(k_0^*, k_1^*)$ . A uniformly random bit is sampled  $b \leftarrow_{\$} \{0, 1\}$ . Queries are generated by  $(\mathbf{q}_{\text{id}}^*)_{\text{ID}} \leftarrow \text{QUERY}(\text{pk}, \text{sk}, k_b^*)$  and  $(\mathbf{q}_{\text{id}}^*)_{\text{ID}}$  is given to  $\mathcal{O}$ .  $\mathcal{O}$  gives  $(\mathbf{q}_{\text{id}}^*)_C$  to  $\mathcal{A}$ .
5.  $\mathcal{A}$  is given more oracle access to QUERY.
6.  $\mathcal{A}$  outputs  $b^* \in \{0, 1\}$ . The output of the experiment is defined to be 1 if  $b^* = b$  and 0 otherwise.

$\Pi^{\text{APIR}}$  is  $(\lambda, t)$ -APIR privacy-preserving for all PPT adversary  $\mathcal{A}$  if there exists a negligible function  $\text{negl}$  such that

$$\Pr [\text{PrivA}_{\mathcal{A}, \Pi^{\text{APIR}}, t}(1^\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda)$$

## 5 Our Scheme: Synchronized APIR

In this section, we present our construction of APIR called *Synchronized APIR*. Synchronized APIR allows the client to synchronize the global “view” of the databases before querying for data, resulting in a high one-time communication cost to initialize the scheme but low amortized computational and communication cost to query. Further, Synchronized APIR makes use of a black-box sPIR scheme described in Section 3.1 in a hybrid PIR, meaning that any improvements to sPIR schemes result in improvements to Synchronized APIR. For the rest of this section, we will introduce the concept of Synchronized APIR in Section 5.1, then follow up with the full description in Section 5.2. Finally, we describe how we implement Synchronized APIR in Section 5.3.

### 5.1 Concept

Let us consider an example of keyword-value databases in Figure 1 step ①, and suppose that each database is independently operated by a PIR server in a network. A client with a given query keyword  $k$  wishes to retrieve the value stored in database  $\text{DB}_0$  associated with keyword  $k$ , while hiding  $k$  from a malicious party controlling the servers. How could this be achieved? In the traditional mPIR setting, databases are required to be consistent for

correctness. However, because this is not the case in Figure 1 step ①, mPIR is not clearly achievable. Instead, the client must resort to the costly sPIR on  $\text{DB}_0$  and disregard  $\text{DB}_1, \text{DB}_2$  and  $\text{DB}_3$ .

To circumvent this issue, Synchronized APIR makes use of some of the keyword-value pairs in  $\text{DB}_1, \text{DB}_2, \text{DB}_3$  that are consistent with  $\text{DB}_0$  to “assist” in reducing the cost of sPIR on  $\text{DB}_0$  via a hybrid PIR. To demonstrate this concept, we will walk through the example in Figure 1 below.

**Step ①:** For the purpose of this demonstration, we assume that the client can fully observe  $\text{DB}_0, \dots, \text{DB}_3$  (how exactly this is done will be explained in the upcoming section). The client first notices that there are keyword-value pairs in  $\text{DB}_1, \text{DB}_2, \text{DB}_3$  inconsistent with  $\text{DB}_0$ , indicated in light grey. The client prefers the “correct” versions of the keyword-value pairs according to  $\text{DB}_0$ , so the inconsistent pairs in  $\text{DB}_1, \text{DB}_2, \text{DB}_3$  are disregarded.

**Step ②:** To lower the cost of Synchronized APIR, the participating parties want to apply traditional mPIR to “assist” wherever possible. This requires that the client first determines how many servers can collude without leaking  $k$ , i.e. the  $t$  collusion threshold similarly to the one in CGKS in Section 3.2. Let suppose that the client decides that  $t = 2$ . Recall from Section 3.2 that a given threshold  $t$  requires at least  $t + 1$  duplicates of the same database values across the databases; so the client must identify the keyword-value pairs with at least 3 duplicates to pass the threshold requirement.

(a, 100), (b, 200), (c, 300), and (d, 400) have at least 3 duplicates ((a, 100) has 4, so the one in  $\text{DB}_3$  is redundant and disregarded). Neither of (e, 500), (f, 600), or (g, 700) meets the requirement, so they are disregarded. The pairs that do not pass the threshold are indicated in dark grey.

**Step ③:** The client splits the pairs that pass the threshold from those that do not. This reveals, on the top rows in the green box, the pairs which can be retrieved with mPIR, and, on the remaining bottom rows in the purple box, the pairs which can only be retrieved with sPIR.

**Step ④:** The client and servers engage in a hybrid PIR protocol. The client makes separate PIR queries for mPIR and sPIR from key  $k$ . All the servers process the mPIR queries, while only the server operating  $\text{DB}_0$  processes the sPIR query.

There are some important details we have omitted here for conceptual simplicity. In the following section, we will expand on the concept of Synchronized APIR to answer these questions:

- How can the client “synchronize” the databases according to step ① - ③ without needing to download them in full, and at a low communication cost?

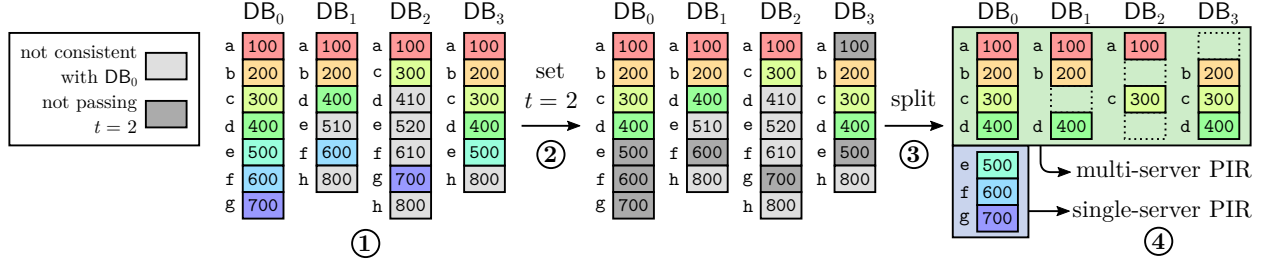


Figure 1: **Concept of Synchronized APIR.** A conceptual demonstration of how Synchronized APIR is able to perform hybrid PIR on inconsistent databases when the collusion threshold is  $t = 2$ .

- How can the client construct a coherent mPIR query in Synchronized APIR when the databases are keyword-value and the mPIR duplicates are scattered across multiple databases?
- How can the communication cost of query and reply be optimized?

## 5.2 Protocol Description

In order to query keyword  $k$ , Client  $\mathcal{C}$  works with main server  $\mathcal{S}_0$  and assisting servers  $\mathcal{S}_1, \dots, \mathcal{S}_n$  through three phases, as illustrated in Figure 2: *Synchronization*, *Setup*, and *Query*. Synchronization and Setup must be completed once at the start, while a new Query session can be repeated for every new query keyword  $\mathcal{C}$  wishes to query.

In the following sections, we will describe step ① - ⑨ of Synchronized APIR as shown in Figure 2 in details.

### 5.2.1 Synchronization Phase

During Synchronization phase,  $\mathcal{C}$  obtains *catalogs* from all servers, and then informs them of the synchronized mPIR keywords that pass the threshold. The goal of this phase is to allow  $\mathcal{C}$  to map out all the keyword-value pairs across the databases without having to download them all in full. We refer to step ① - ③ in Figure 2 and ① - ② in 3 as a demonstration to describe Synchronization phase below.

#### Notations

- $\text{Map}(\cdot)$ , a keyword-value map : suppose that  $M = \text{Map}(S)$ , then  $(k, v) \in S$  iff  $M[k] = v$ .
- $\phi$ , an empty set

#### Input

- ▷ DB<sub>id</sub> **from**  $\mathcal{S}_{id}, \forall id \in \text{ID}$
- ▷  $t$  **from**  $\mathcal{C}$

#### Output

- ▷  $\text{Cat}_{id}, M_{id}$  **to**  $\mathcal{S}_{id}, \forall id \in \text{ID}$
- ▷  $\text{Cat}_0, (\text{Cat}'_{id})_{\text{AID}}, (M_{id})_{\text{ID}}, S$  **to**  $\mathcal{C}$

### Synchronization Protocol

- ①  $\mathcal{C}$  obtains the full catalog  $\text{Cat}_0$  from  $\mathcal{S}_0$  and catalog intersections  $\text{Cat}'_1, \dots, \text{Cat}'_n$  from  $\mathcal{S}_1, \dots, \mathcal{S}_n$  by following the steps below:

- A. For each  $id \in \text{ID}$ ,  $\mathcal{S}_{id}$  generates catalog

$$\text{Cat}_{id} \leftarrow \text{Map}\{(k, h) \mid (k, v) \in \text{DB}_{id}, h = H(v)\}$$

where  $H(\cdot)$  is a universal hash function with short hash values, chosen at random by  $\mathcal{C}$  from the universal family. (A shorter hash length will reduce the communication cost but increases the probability of collision.) This results in the catalogs that are smaller than the databases in size but capable of representing the uniqueness of database values. In Figure 3, we choose  $H(v) = v/10$  to demonstrate the point that an appropriate hash function should be able to produce short unique hash values. In general when the patterns of database values are not predetermined, universal hash functions should be applied to reduce hash collisions.

- B.  $\mathcal{C}$  downloads  $\text{Cat}_0$  from  $\mathcal{S}_0$ . For each  $id \in \text{AID}$ ,  $\mathcal{C}$  engages in catalog intersection protocol  $\text{CATINTERSECTION}$  with  $\mathcal{S}_{id}$ , where  $\mathcal{C}$  obtains catalog intersections  $\text{Cat}'_{id} = \text{Cat}_{id} \cap \text{Cat}_0$  at the end. The goal of  $\text{CATINTERSECTION}$  is to transmit catalog intersections at a cost lower than sending full catalogs.  $\text{CATINTERSECTION}$  is described in Section 5.2.5. Figure 3 demonstrates how keyword-value pairs inconsistent with  $\text{DB}_0$  in  $\text{DB}_1, \dots, \text{DB}_n$  and with  $\text{Cat}_0$  in  $\text{Cat}_1, \dots, \text{Cat}_n$  are eliminated via catalog intersection. With the full catalog and catalog intersections,  $\mathcal{C}$  now has complete information of the keyword-value pairs consistent with  $\text{DB}_0$ .

- ②  $\mathcal{C}$  tags the catalogs to categorize keyword-value pairs for either mPIR or sPIR:

$$((M_{id})_{\text{ID}}, S) \leftarrow \text{TAG}(\text{Cat}_0, (\text{Cat}'_{id})_{\text{AID}}, t)$$

where  $t$  is the collusion threshold i.e. the number of colluding servers up to which  $\mathcal{C}$  can tolerate without

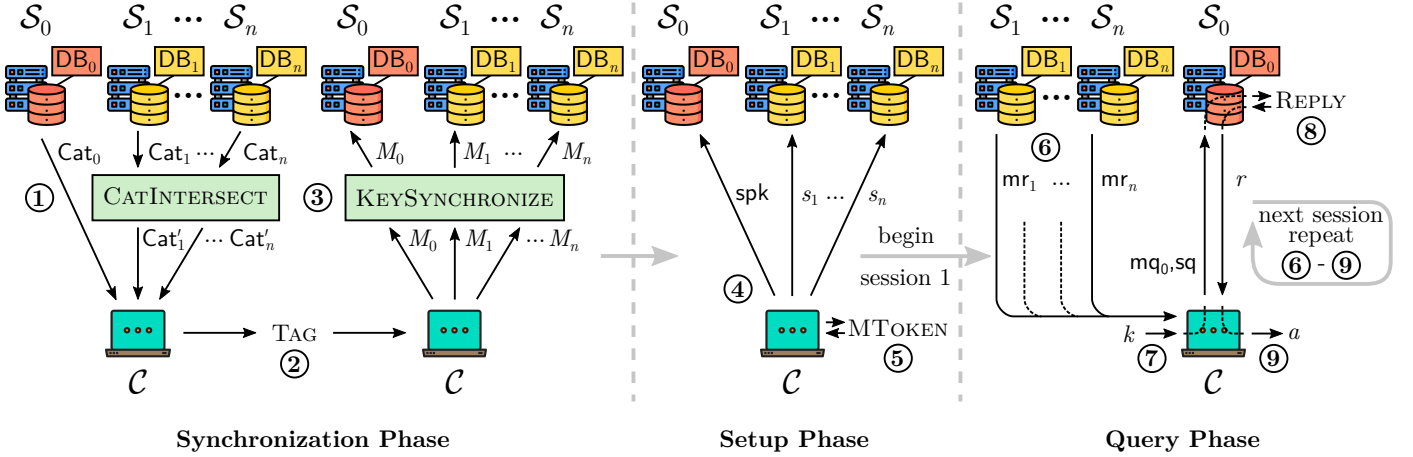


Figure 2: **Overview of Synchronized APIR.** Synchronized APIR comprises three distinct phases between client  $\mathcal{C}$ , main server  $\mathcal{S}_0$ , and assisting Servers  $\mathcal{S}_1 - \mathcal{S}_n$ : Synchronization, Setup, and Query. Synchronization and Setup are completed once at the start, while a new Query session can be repeated with every new query keyword  $k$ .

leaking query keywords;  $M_{id}$  is the mPIR keyword set i.e. set of keywords determined to be queried via mPIR from  $\mathcal{S}_{id}$ ;  $S$  is the sPIR keyword set i.e. set of keywords determined to be queried via sPIR from  $\mathcal{S}_0$ . TAG is detailed in Algorithm 1 below. In Figure 3,  $\mathcal{C}$  has decided that  $t = 2$ , so the pairs that pass the threshold for mPIR must have  $t + 1 = 3$  duplicates: 1 duplicate in  $Cat_0$  and 2 duplicates in  $Cat'_1, Cat'_2, Cat'_3$ . (a, 10), (b, 20), (c, 30), and (d, 40) all pass the threshold ((a, 10) in  $Cat'_3$  is redundant and disregarded). Neither of (e, 50), (f, 60), or (g, 70) pass the threshold, so  $\mathcal{C}$  disregards them. This results in  $M_0, \dots, M_4$  for the set of keywords that pass the threshold for mPIR, and the rest in  $S$  for sPIR.

**Algorithm 1 Catalog Tagging.** In line 5 - 8, if there are at least  $t$  duplicates of a given keyword, then  $t$  of them are tagged for mPIR and stored in  $M_{id}$  ( $t + 1$  in total, including those in  $M_0$ ). In line 10, the rest of the keywords are tagged for sPIR. In line 6,  $C_k$  can be chosen at random or with a specific optimization strategy.

```

1: procedure TAG( $Cat_0, (Cat'_{id})_{AID}, t$ )
2:    $\forall id \in ID : M_{id} \leftarrow \phi$ 
3:   for  $k \in Keys(Cat_0)$  do
4:      $G \leftarrow \{id \in AID \mid k \in Keys(Cat'_{id})\}$ 
5:     if  $|G| \geq t$  then
6:       choose  $C_k \subseteq G$  such that  $|C_k| = t$ 
7:        $\forall id \in C_k \cup \{0\} : M_{id} \leftarrow M_{id} \cup \{k\}$ 
8:     end if
9:   end for
10:   $S \leftarrow Keys(Cat_0) \setminus M_0$ 
11:  return  $((M_{id})_{ID}, S)$ 
12: end procedure

```

- ③ For each  $id \in ID$ ,  $\mathcal{C}$  and  $\mathcal{S}_{id}$  engage in protocol KEYSYNCHRONIZE, where  $\mathcal{S}_{id}$  obtains  $M_{id}$  at the end. The goal of KEYSYNCHRONIZE is to transmit  $M_{id}$  at a cost lower than sending them in full. KEYSYNCHRONIZE is described in Section 5.2.6.

## 5.2.2 Setup Phase

Let  $l = \text{poly}(\lambda)$  denote the total number of queries that  $\mathcal{C}$  will be making during Query phase ( $l$  is not necessarily predetermined at this point). We call each query during Query phase a Query *session*, each denoted with session ID  $sid \in \{1, \dots, l\}$ .

### Notations

- $sid \in \{1 \dots l\}$ , session ID. Starting at  $sid = 1$ ,  $sid$  increases by 1 for each new iteration of Query. Let  $l$  denote the number of sessions and  $[l] := \{1 \dots l\}$  the set of all session ID's.
- $\text{index}(a, A) := |\{b \in A \mid b < a\}|$ , index of  $a$  in set  $A$ .

During Setup,  $\mathcal{C}$  sends each of  $\mathcal{S}_1, \dots, \mathcal{S}_n$  a random PRG seed which will be used to generate random mPIR queries for all Query sessions  $sid \in [l]$ . We follow step ④ and ⑤ in Figure 2 and ⑤ in 4 to describe Setup phase below.

### Input

- ▷  $|M_{id}|$  **from**  $\mathcal{S}_{id}, \forall id \in AID$
- ▷  $l, (M_{id})_{AID}$  **from**  $\mathcal{C}$

### Output

- ▷  $spk$  **to**  $\mathcal{S}_0$
- ▷  $(mq_{id}^{sid})_{sid \in [l]}$  **to**  $\mathcal{S}_{id}, \forall id \in AID$

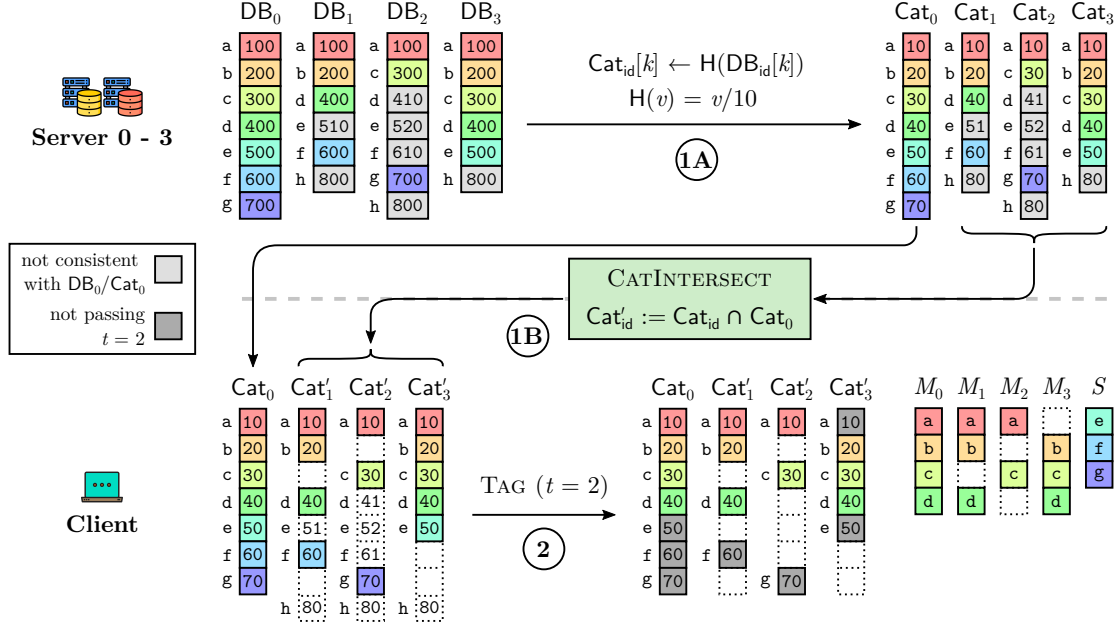


Figure 3: **Synchronization.** An example of databases in a Synchronized APIR setting with 3 assisting servers that undergo Synchronization phase.

▷  $\text{spk}, \text{ssk}, (\text{tok}^{\text{sid}})_{[l]}$  to  $\mathcal{C}$

### Setup Protocol

- ④  $\mathcal{C}$  generates an sPIR key pair  $(\text{spk}, \text{ssk}) \leftarrow \text{SGEN}(\lambda)$  and sends  $\text{spk}$  to main server  $\mathcal{S}_0$ .

For each assisting servers  $\text{id} \in \text{AID}$ ,  $\mathcal{C}$  samples a PRG seed  $s_{\text{id}} \in \{0, 1\}^\lambda$  uniformly at random and sends  $s_{\text{id}}$  to  $\mathcal{S}_{\text{id}}$ .  $\mathcal{C}$  and each  $\mathcal{S}_{\text{id}}$  define

$$(\text{mq}_{\text{id}}^{\text{sid}})_{\text{sid} \in [l]} := \text{PRG}(s_{\text{id}}, l \cdot |M_{\text{id}}|)$$

where  $l \cdot |M_{\text{id}}|$  is the total length of the output string and  $|\text{mq}_{\text{id}}^{\text{sid}}| = |M_{\text{id}}|$ .

Note that in practice when the PRG is implemented with a stream cipher,  $l$  does not need to be predetermined and each  $\text{mq}_{\text{id}}^{\text{sid}}$  can be generated on-the-fly.

- ⑤  $\mathcal{C}$  defines query *tokens*

$$\text{tok}^{\text{sid}} := \text{MTOKEN}((M_{\text{id}})_{\text{ID}}, (\text{mq}_{\text{id}}^{\text{sid}})_{\text{AID}})$$

for all  $\text{sid} \in [l]$ . MTOKEN is defined in Algorithm 2 below. Intuitively, a token is an XOR of the random queries  $\text{mq}_{\text{id}}^{\text{sid}}$  for all  $\text{id} \in \text{AID}$  when the bits are aligned to the ordering of the corresponding keywords in  $M_0$ . This is visually demonstrated in Figure 4 step ⑤.

Likewise as in the previous step, each  $\text{tok}^{\text{sid}}$  can be generated on-the-fly.

### Algorithm 2 mPIR Token Generation.

- 1: **procedure** MTOKEN( $((M_{\text{id}})_{\text{ID}}, (\text{mq}_{\text{id}}^{\text{sid}})_{\text{AID}})$ )
- 2:  $\text{tok} \leftarrow \{0\}^{|M_0|}$
- 3: **for**  $k \in M_0$  **do**
- 4:  $C_k \leftarrow \{\text{id} \in \text{AID} : k \in M_{\text{id}}\}$
- 5:  $\forall \text{id} \in C_k \cup \{0\} : i_{\text{id}} \leftarrow \text{index}(k, M_{\text{id}})$
- 6:  $\text{tok}[i_0] \leftarrow \bigoplus_{\text{id} \in C} \text{mq}_{\text{id}}[i_{\text{id}}]$
- 7: **end for**
- 8: **return** tok
- 9: **end procedure**

### 5.2.3 Special Case: mPIR-Only Query Phase

Before proceeding to describing the full protocol of Query phase, it would be instructive to walk through the special case in which all keyword-value pairs can be retrieved with mPIR i.e.  $S = \phi$  in Algorithm 1. The goal is to show how all the mPIR components fit together during Query phase to provide correct answers, without having to deal with the complexity of the hybrid PIR. We shall follow the steps in Figure 4 for this walk-through.

Suppose that at this point in time  $\mathcal{C}$  and  $\mathcal{S}_0, \dots, \mathcal{S}_n$  have already completed Synchronization in Section 5.2.1 and Setup in Section 5.2.2. We will now pick up step ⑥ and ⑦\*-⑨\* (\* to denote the steps for this special case) in Figure 4 from here.

#### Input

- ▷  $\text{DB}_{\text{id}}, M_{\text{id}}$  from  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{ID}$   
 ▷  $\text{mq}_{\text{id}}^{\text{sid}}$  from  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{AID}$



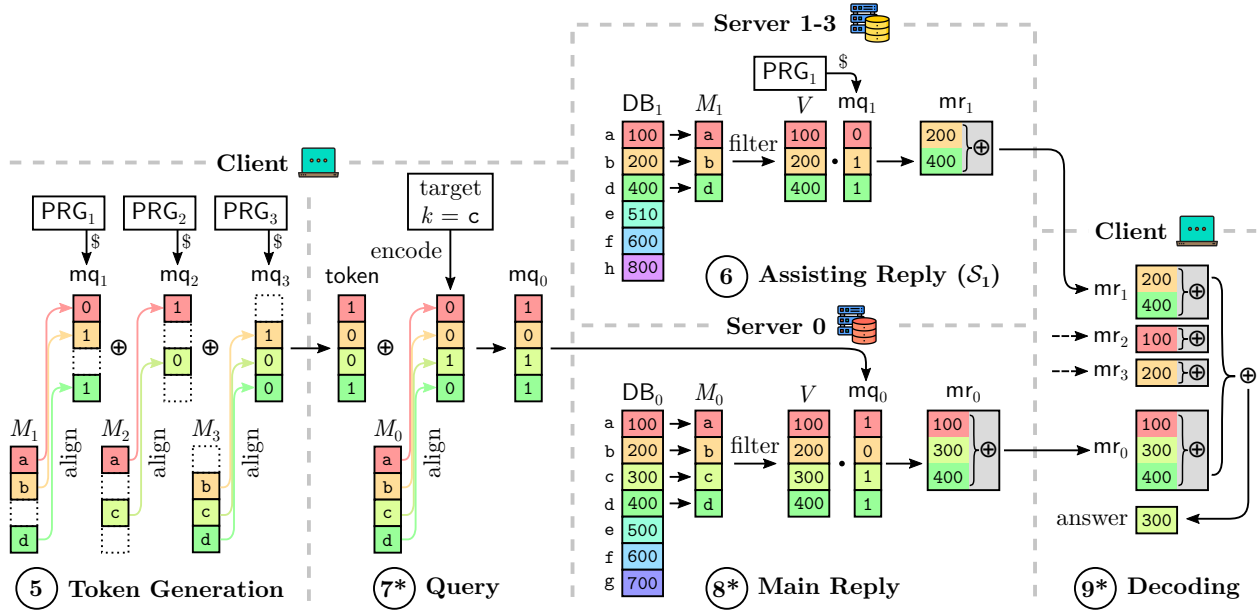


Figure 4: **mPIR-Only Query Phase Workflow.** Workflow of mPIR components in during Query phase to demonstrate how mPIR keyword-value pairs can be queried with the client’s query keyword.

▷  $M_0, \text{tok}^{\text{sid}}, k^{\text{sid}}$  from  $\mathcal{C}$

**Output**

▷  $a^{\text{sid}}$  to  $\mathcal{C}$

### mPIR-Only Query Protocol for Session $\text{sid}$

⑥ For each assisting servers  $\text{id} \in \text{AID}$ :

A.  $\mathcal{S}_{\text{id}}$  generates a reply from the query

$$\text{mr}_{\text{id}}^{\text{sid}} \leftarrow \text{MREPLY}(M_{\text{id}}, \text{DB}_{\text{id}}, \text{mq}_{\text{id}}^{\text{sid}})$$

and sends  $\text{mr}_{\text{id}}^{\text{sid}}$  to  $\mathcal{C}$ . MREPLY is described in Algorithm 3 below. Figure 4 step ⑥ demonstrates the mPIR reply generation process: the values of  $\text{DB}_{\text{id}}$  are filtered with  $M_{\text{id}}$  by key; and the dot product in  $\text{GF}(2)$  between the filtered values and random query bits produces the reply. This process is identical for  $\mathcal{S}_2$  and  $\mathcal{S}_3$ .

**Algorithm 3 mPIR Reply Generation.** The dot product in line 3 is defined in  $\text{GF}(2)$ .

```

1: procedure MREPLY( $M_{\text{id}}, \text{DB}_{\text{id}}, \text{mq}_{\text{id}}$ )
2:    $V \leftarrow (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}}$ 
3:    $\text{mr}_{\text{id}} \leftarrow \text{mq}_{\text{id}} \cdot V$ 
4:   return  $\text{mr}_{\text{id}}$ 
5: end procedure

```

Note that this step can be completed offline as it does not require a query keyword  $k^{\text{sid}}$  by  $\mathcal{C}$ . This means

that during online time (that is, step ⑦\* onwards),  $\mathcal{C}$  only needs to interact with  $\mathcal{S}_0$ , which significantly reduces the latency.

⑦\* Once  $\mathcal{C}$  has decided on a mPIR query keyword  $k^{\text{sid}} \in M_0$  to query, she generates a mPIR query for  $\mathcal{S}_0$

$$\text{mq}_0^{\text{sid}} \leftarrow \text{MQQUERY}(M_0, \text{tok}^{\text{sid}}, k^{\text{sid}})$$

where MQQUERY is described in Algorithm 4 below. Intuitively, as visually demonstrated in Figure 4, MQQUERY flips one bit of the token  $\text{tok}^{\text{sid}}$  at the index position of  $k^{\text{sid}}$  in  $M_0$  to produce the query  $\text{mq}_0^{\text{sid}}$ .

Finally,  $\mathcal{C}$  sends  $\text{mq}_0^{\text{sid}}$  to  $\mathcal{S}_0$  to query.

### Algorithm 4 mPIR Query Generation.

```

1: procedure MQQUERY( $M_0, \text{tok}, k$ )
2:    $\text{mq}_0 \leftarrow \text{tok}$ 
3:   if  $k \in M_0$  then
4:      $i \leftarrow \text{index}(k, M_0)$ 
5:      $\text{mq}_0[i] \leftarrow \text{mq}_0[i] \oplus 1$ 
6:   end if
7:   return  $\text{mq}_0$ 
8: end procedure

```

⑧\* Upon receiving  $\text{mq}_0^{\text{sid}}$ ,  $\mathcal{S}_0$  processes the query in the same way queries are processed in step ⑥, except that here  $\text{mq}_0^{\text{sid}}$  is given by  $\mathcal{C}$ ,

$$\text{mr}_0^{\text{sid}} \leftarrow \text{MREPLY}(\text{mq}_0^{\text{sid}}, M_0, \text{DB}_0)$$

$\mathcal{S}_0$  returns  $\text{mr}_0^{\text{sid}}$  to  $\mathcal{C}$ .

- ⑨\* Now that  $\mathcal{C}$  has received all  $\text{mr}_0^{\text{sid}}, \dots, \text{mr}_n^{\text{sid}}$ , she can decode them for the answer by a simple XOR

$$a^{\text{sid}} \leftarrow \bigoplus_{\text{id} \in \text{ID}} \text{mr}_{\text{id}}^{\text{sid}}$$

where  $a^{\text{sid}} = \text{DB}[k^{\text{sid}}]$  as a result.

To see why it is the case that  $a^{\text{sid}} = \text{DB}[k^{\text{sid}}]$ , recall from step ⑤ and ⑦\* in Figure 4. If we “peel off”  $\text{mq}_1^{\text{sid}}, \dots, \text{mq}_n^{\text{sid}}$  from  $\text{mq}_0^{\text{sid}}$  via XOR, the result is exactly the encoding of  $k^{\text{sid}}$  by design. Because all keyword-value pairs of  $\text{DB}_1, \dots, \text{DB}_n$  are consistent with that of  $\text{DB}_0$ , in step ⑨\* all other non-target values are peeled off except for the target value.

The privacy of the  $k^{\text{sid}}$  is intact if no more than  $t$  Servers collude. This is because up to  $t$  Servers can observe up to  $t$  of  $\text{mq}_{\text{id}}^{\text{sid}}$ 's which are pseudorandom as intended by the collusion threshold during TAG.

#### 5.2.4 Query Phase

In this section, we will expand on the special case in the last section to describe the full Query phase when  $S \neq \phi$  in Algorithm 1. This enables the client to retrieve any values in  $\text{DB}_0$  without leaking to the servers whether the query keyword is in sPIR or mPIR. Below, we follow Figure 5 to describe step ⑥ - ⑨. We refer to step ⑥, ⑦\* - ⑨\* from Section 5.2.3.

##### Input

- ▷  $\text{DB}_0, M_0, \text{spk}$  **from**  $\mathcal{S}_0$
- ▷  $\text{DB}_{\text{id}}, M_{\text{id}}, \text{mq}_{\text{id}}^{\text{sid}}$  **from**  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{AID}$
- ▷  $M_0, S, \text{spk}, \text{ssk}, \text{tok}^{\text{sid}}, k^{\text{sid}}$  **from**  $\mathcal{C}$

##### Output

- ▷  $a^{\text{sid}}$  **to**  $\mathcal{C}$

#### Query Protocol for Session sid

- ⑥ This is the same as step ⑥ in Section 5.2.3.
- ⑦  $\mathcal{C}$  decides on a query keyword  $k^{\text{sid}} \in M_0 \cup S$  and generates a hybrid-PIR query for  $\mathcal{S}_0$

$$(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}}) \leftarrow \text{QUERY}(M_0, S, \text{spk}, \text{tok}^{\text{sid}}, k^{\text{sid}})$$

where  $\text{sq}^{\text{sid}}$  is the sPIR query and  $\text{mq}_0^{\text{sid}}$  mPIR query. QUERY is described in Algorithm 5 below.  $\mathcal{C}$  sends  $(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})$  to  $\mathcal{S}_0$ .

The hybridization of the mPIR-sPIR query is visually demonstrated in Figure 5. First, step ⑦\* is followed to generate  $\text{mq}_0^{\text{sid}}$ . (Note that if  $k^{\text{sid}} \notin M_0$ , then  $\text{mq}_0^{\text{sid}}$  is a “blank” token as per Algorithm 4.)

Next, if  $k^{\text{sid}} \in S$ , then  $k^{\text{sid}}$  is encoded by its index position in  $S$  as an input to SQUERY; otherwise if  $k^{\text{sid}} \in M_0$ , then it is encoded as the last index position as an input to SQUERY. In this example, because  $k^{\text{sid}} = c \in M_0$ , it is encoded as 3, the 0-based index of the last position.

#### Notations

- $\perp$ , an empty string/value

#### Algorithm 5 Query Generation.

```

1: procedure QUERY( $M_0, S, \text{spk}, \text{tok}, k$ )
2:    $\text{mq}_0 \leftarrow \text{MQUERY}(M_0, \text{tok}, k)$ 
3:    $\text{sq} \leftarrow \perp$ 
4:   if  $S \neq \phi$  then
5:     if  $k \in S$  then
6:        $i \leftarrow \text{index}(k, S)$ 
7:        $\text{sq} \leftarrow \text{SQUERY}(\text{spk}, i, |S| + 1)$ 
8:     else
9:        $\text{sq} \leftarrow \text{SQUERY}(\text{spk}, |S|, |S| + 1)$ 
10:    end if
11:  end if
12:  return  $(\text{mq}_0, \text{sq})$ 
13: end procedure

```

- ⑧  $\mathcal{S}_0$  receives the query  $(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})$  from  $\mathcal{C}$  and generates the reply

$$r^{\text{sid}} \leftarrow \text{REPLY}(M_0, \text{DB}_0, \text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})$$

where  $r^{\text{sid}}$  is the resulting hybrid PIR reply. REPLY is described in Algorithm 6 below.  $\mathcal{S}_0$  returns  $r^{\text{sid}}$  to  $\mathcal{C}$ .

#### Algorithm 6 Reply Generation.

```

1: procedure REPLY( $M_0, \text{DB}_0, \text{mq}_0, \text{sq}$ )
2:    $\text{mr}_0 \leftarrow \text{MREPLY}(M_0, \text{DB}_0, \text{mq}_0)$ 
3:    $S \leftarrow \text{Keys}(\text{DB}_0) \setminus M_0$ 
4:    $r \leftarrow \perp$ 
5:   if  $S \neq \phi$  then
6:      $V \leftarrow (\text{DB}_0[k])_{k \in S}$ 
7:      $r \leftarrow \text{SREPLY}(\text{spk}, V \parallel \text{mr}_0, \text{sq})$ 
8:   else
9:      $r \leftarrow \text{mr}_0$ 
10:  end if
11:  return  $r$ 
12: end procedure

```

Figure 5 demonstrates the hybridized reply generation process. First, the mPIR reply is generated according to step ⑧\*, resulting in  $\text{mr}_0^{\text{sid}}$ . (Note that

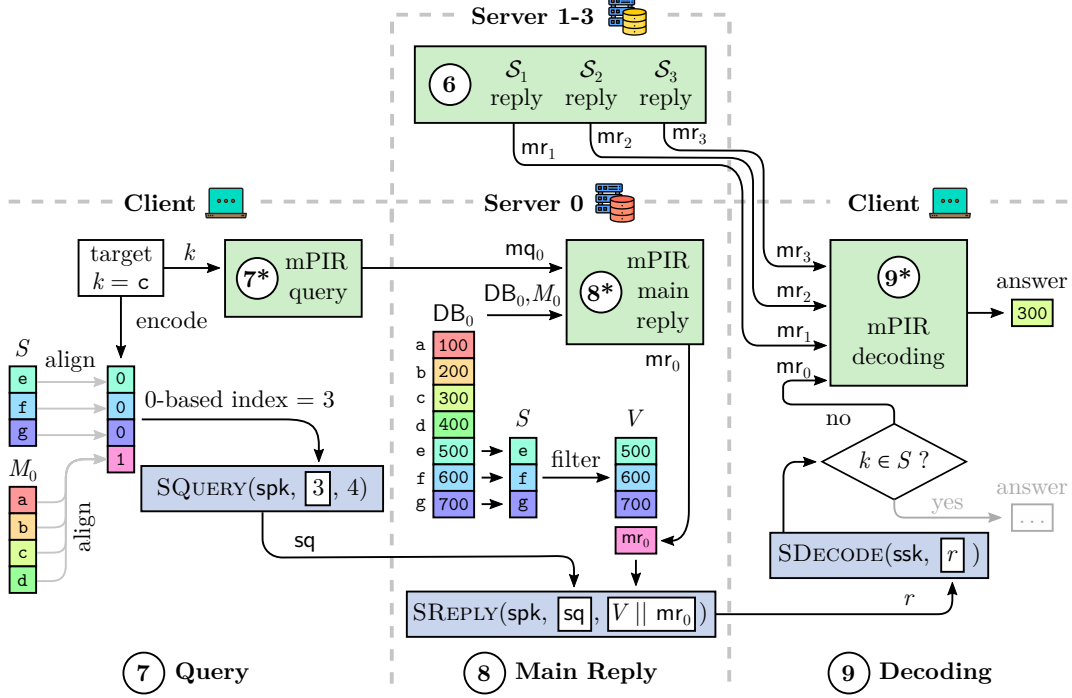


Figure 5: **Query Phase Workflow** Workflow of hybrid components in during Query phase to demonstrate how keyword-value pairs in  $DB_0$  can be queried with the client’s query key.

if  $mq_0^{\text{sid}}$  is not a blank query, then  $mr_0^{\text{sid}}$  is a genuine mPIR reply; otherwise, it is a random value. This fact, however, is unknown to  $\mathcal{S}_0$ .) Next, the values of  $DB_0$  is filtered with  $S$  by key, resulting in  $V$ .  $mr_0^{\text{sid}}$  is appended to  $V$  as the last item, i.e.  $V || mr_0^{\text{sid}}$ . The sPIR query  $sq^{\text{sid}}$  and  $V || mr_0^{\text{sid}}$  are input to the sPIR reply generation algorithm  $SREPLY$  to produce a reply.

At this point, it becomes clear why the sPIR query  $sq$  is generated the way it is as described in step ⑦. When  $k^{\text{sid}} \in S$ ,  $sq_{\text{sid}}$  is generated to target values in  $V$ , ignoring  $mr_0^{\text{sid}}$  which is the last item. However, when  $k \in M_0$ ,  $sq_{\text{sid}}$  is generated to target the last value, which is where  $mr_0^{\text{sid}}$  is located. The resulting reply  $r^{\text{sid}}$  therefore encodes either a targeted value in  $V$  when  $k^{\text{sid}} \in S$  or  $mr_0^{\text{sid}}$  when  $k^{\text{sid}} \in M_0$ .

- ⑨  $\mathcal{C}$  receives the reply  $r^{\text{sid}}$  from  $\mathcal{S}_0$  and decodes it using the same query keyword  $k^{\text{sid}}$  from step ⑦ to obtain the answer

$$a^{\text{sid}} \leftarrow \text{DECODE} \left( S, \text{ssk}, k^{\text{sid}}, (mr_{\text{id}}^{\text{sid}})_{\text{id} \in \text{AID}}, r^{\text{sid}} \right)$$

DECODE is described in Algorithm 7 below.

Figure 5 demonstrates how all the replies from  $\mathcal{S}_0, \dots, \mathcal{S}_3$  come together for  $\mathcal{C}$  to decode the answer.

#### Algorithm 7 Decoding.

```

1: procedure DECODE( $S, \text{ssk}, k, (mr_{\text{id}}^{\text{sid}})_{\text{id} \in \text{AID}}, r$ )
2:   if  $k \in S$  then
3:      $a \leftarrow \text{SDECODE}(\text{ssk}, r)$ 
4:   else
5:     if  $S \neq \phi$  then
6:        $mr_0 \leftarrow \text{SDECODE}(\text{ssk}, r)$ 
7:     else
8:        $mr_0 \leftarrow r$ 
9:     end if
10:     $a \leftarrow \bigoplus_{\text{id} \in \text{ID}} mr_{\text{id}}$ 
11:  end if
12:  return  $a$ 
13: end procedure

```

First, the reply  $r^{\text{sid}}$  from  $\mathcal{S}_0$  is decoded using secret key  $\text{ssk}$  in the sPIR decoding algorithm  $SDECODE$ . If  $k^{\text{sid}} \in S$ , then  $\mathcal{C}$  already has the final answer; otherwise,  $\mathcal{C}$  obtains  $mr_0^{\text{sid}}$  and follows step ⑨\*. Here,  $\mathcal{C}$  XOR all  $mr_0^{\text{sid}}, \dots, mr_3^{\text{sid}}$  to obtain the answer.

For privacy, it is important that regardless of whether  $k^{\text{sid}} \in S$ ,  $mr_1^{\text{sid}}, \dots, mr_n^{\text{sid}}$  must be downloaded by  $\mathcal{C}$  from  $\mathcal{S}_1, \dots, \mathcal{S}_n$ , respectively. If this step is skipped, then  $\mathcal{S}_1, \dots, \mathcal{S}_n$  can infer that  $k^{\text{sid}} \in S$ .

### 5.2.5 Catalog Intersection

Catalog Intersection protocol is a part of Synchronization phase to reduce the communication cost of transferring catalog intersections. The idea is as follows: given that  $\mathcal{C}$  has already obtained  $\text{Cat}_0$ , then  $\text{Cat}'_{\text{id}} = \text{Cat}_{\text{id}} \cap \text{Cat}_0$  can be compressed with a hash function (we call these hashes *digests*) for transporting. We describe Catalog Intersection protocol below.

#### Input

- ▷  $\text{Cat}_{\text{id}}$  **from**  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{AID}$
- ▷  $\text{Cat}_0$  **from**  $\mathcal{C}$

#### Output

- ▷  $(\text{Cat}'_{\text{id}})_{\text{AID}}, (\text{Dig}_{\text{id}})_{\text{AID}}$  **to**  $\mathcal{C}$

### Catalog Intersection Protocol

- ① For each  $\text{id} \in \text{AID}$ ,  $\mathcal{S}_{\text{id}}$  generates digests from  $\text{Cat}_{\text{id}}$

$$\text{Dig}_{\text{id}} \leftarrow \{G((k, h)) \mid (k, h) \in \text{Cat}_{\text{id}}\}$$

where  $G(\cdot)$  is a universal hash function with short hash values, chosen at random from the universal family by  $\mathcal{C}$ . (A shorter hash length will reduce the communication cost but increase the probability of collision.) A digest is therefore a hashed representation of a keyword-hash pair in the catalogs.

$\mathcal{S}_{\text{id}}$  sends the digest set  $\text{Dig}_{\text{id}}$  to  $\mathcal{C}$ .

- ②  $\mathcal{C}$  downloads all the digest sets  $(\text{Dig}_{\text{id}})_{\text{AID}}$  from  $\mathcal{S}_1, \dots, \mathcal{S}_n$ . For each  $\text{id} \in \text{AID}$ ,  $\mathcal{C}$  remaps digest sets into catalog intersections as follows:

$$\text{Cat}'_{\text{id}} \leftarrow \text{Map}\{(k, h) \in \text{Cat}_0 \mid G((k, h)) \in \text{Dig}_{\text{id}}\}$$

### 5.2.6 Keyword Synchronization

Keyword Synchronization protocol is a part of Synchronization phase to reduce the communication cost of transferring mPIR keyword sets  $M_{\text{id}}$  from  $\mathcal{C}$  to  $\mathcal{S}_{\text{id}}$ . The idea is to map each keyword in  $M_{\text{id}}$  to its corresponding index in  $\text{Dig}_{\text{id}}$  for transport, which can then be remapped to  $M_{\text{id}}$  on the server side. We describe Keyword Synchronization protocol below.

#### Input

- ▷  $\text{Dig}_{\text{id}}, \text{Cat}_{\text{id}}$  **from**  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{ID}$
- ▷  $(M_{\text{id}})_{\text{ID}}, (\text{Dig}_{\text{id}})_{\text{ID}}$  **from**  $\mathcal{C}$

#### Output

- ▷  $M_{\text{id}}$  **to**  $\mathcal{S}_{\text{id}}, \forall \text{id} \in \text{ID}$

### Keyword Synchronization Protocol

- ① For all  $\text{id} \in \text{ID}$ ,  $\mathcal{C}$  converts  $M_{\text{id}}$  to the set of indices corresponding to the ordering of  $\text{Keys}(\text{Cat}_0)$  for main server and of  $\text{Dig}_{\text{id}}$  for assisting servers:

$$I_0 \leftarrow \text{KEYTOINDEX}_0(M_0, \text{Cat}_0)$$

$$\forall \text{id} \in \text{AID} : I_{\text{id}} \leftarrow \text{KEYTOINDEX}_{\text{id}}(M_{\text{id}}, \text{Dig}_{\text{id}}, \text{Cat}_{\text{id}})$$

$\text{KEYTOINDEX}$  is described in Algorithm 8 below.  $\mathcal{C}$  then sends  $I_{\text{id}}$  to  $\mathcal{S}_{\text{id}}$ .  $I_{\text{id}}$  is a small, indexed representation of  $M_{\text{id}}$  which can be transported at a low communication cost.

---

#### Algorithm 8 Key-to-Index Mapping

---

- 1: **procedure**  $\text{KEYTOINDEX}_0(M_0, \text{Cat}_0)$
  - 2:  $I_0 \leftarrow \{i \mid k \in M_0, i = \text{index}(k, \text{Keys}(\text{Cat}_0))\}$
  - 3: **return**  $I_0$
  - 4: **end procedure**
  - 5: **procedure**  $\text{KEYTOINDEX}_{\text{id}}(M_{\text{id}}, \text{Dig}_{\text{id}}, \text{Cat}_{\text{id}})$
  - 6:  $D \leftarrow \{G((k, h)) \mid k \in M_{\text{id}}, h = \text{Cat}_{\text{id}}[k]\}$
  - 7:  $I_{\text{id}} \leftarrow \{i \mid d \in D, i = \text{index}(d, \text{Dig}_{\text{id}})\}$
  - 8: **return**  $I_{\text{id}}$
  - 9: **end procedure**
- 

- ② For each  $\text{id} \in \text{ID}$ ,  $\mathcal{S}_{\text{id}}$  receives  $I_{\text{id}}$  and remaps it back to  $M_{\text{id}}$

$$M_0 \leftarrow \text{INDEXTOKEY}_0(I_0, \text{Cat}_0)$$

$$\forall \text{id} \in \text{AID} : M_{\text{id}} \leftarrow \text{INDEXTOKEY}_{\text{id}}(I_{\text{id}}, \text{Dig}_{\text{id}}, \text{Cat}_{\text{id}})$$

$\text{INDEXTOKEY}$  is described in Algorithm 9 below.

---

#### Algorithm 9 Index-to-Key Mapping

---

- 1: **procedure**  $\text{INDEXTOKEY}_0(I_0, \text{Cat}_0)$
  - 2:  $K \leftarrow \text{Keys}(\text{Cat}_0)$
  - 3:  $M_0 \leftarrow \{k \in K \mid i \in I_0, \text{index}(k, K) = i\}$
  - 4: **return**  $M_0$
  - 5: **end procedure**
  - 6: **procedure**  $\text{INDEXTOKEY}_{\text{id}}(I_{\text{id}}, \text{Dig}_{\text{id}}, \text{Cat}_{\text{id}})$
  - 7:  $D \leftarrow \{d \in \text{Dig}_{\text{id}} \mid i \in I_{\text{id}}, \text{index}(d, \text{Dig}_{\text{id}}) = i\}$
  - 8:  $M_{\text{id}} \leftarrow \{k \mid (k, h) \in \text{Cat}_{\text{id}}, G((k, h)) \in D\}$
  - 9: **return**  $M_{\text{id}}$
  - 10: **end procedure**
- 

### 5.2.7 Catalog Intersection and Keyword Synchronization in Relation to Set Reconciliation

Set reconciliation is an active research area to find the symmetric difference between two sets from two remote parties over a network at a low communication cost. Notably, Eppstein *et al.* [8] propose a practical scheme whose

communication cost is asymptotically linear in the size of the symmetric difference. One therefore could wonder if these schemes would apply to our Catalog Intersection and Keyword Synchronization problem. This is because Catalog Intersection and Keyword Synchronization are essentially about finding remote set intersections, which is implied by set symmetric differences. We argue that our solution as presented above solve this problem more efficiently by killing two birds with one stone: instead of apply two separate instances of set reconciliation to Catalog Intersection and Keyword Synchronization (and therefore paying the cost twice), we use the ordering information provided by digest sets  $\text{Dig}_{\text{id}}$  to reduce keyword sets  $M_{\text{id}}$  to sets of indices, and therefore reduce the upload cost by the client (which is especially important in an asymmetric network connection). It is not clear how this cost-saving technique can be applied to existing set reconciliation schemes in literature.

### 5.2.8 Keyword Compression

When query keywords are long, this can result a high communication cost during Synchronization phase since all keywords in catalogs must be transported to the client. For this reason, we aim to reduce this cost by representing each query keyword  $k$  by its hash  $\tilde{H}(k)$ , where  $\tilde{H}(\cdot)$  is a universal hash function with short hash values, chosen from the universal family by the client. When querying keyword  $k$ , the client computes the hash and uses it as a keyword to query instead.

Despite the cost reduction, this can introduce a new problem when the keyword space is large, which can either lead to long hash values or high collision (i.e. high false-positive) rate. Our suggested workaround for this problem is the keep to hash values short (and therefore high collusion rate) but let the servers modify the databases by prepending the keyword to the value so that a false positive can be detected at the end. That is,

$$\widetilde{\text{DB}}_{\text{id}} := \text{Map} \{ (k, \tilde{v}) \mid (k, v) \in \text{DB}_{\text{id}}, \tilde{v} := k \| v \}$$

When the client receives the answer  $\tilde{v}$  at the end, she can check whether the obtained keyword from the answer is the same as the query keyword. This, indeed, comes at the cost of the servers processing extra loads and longer database values.

## 5.3 Implementation

Synchronized APIR is implemented in Rust, where keyword compression in Section 5.2.8 is implemented by default. The universal hash functions  $\text{H}$ ,  $\text{G}$ , and  $\tilde{\text{H}}$  are instantiated with Google’s CityHash (<https://github.com/google/cityhash>), each seeded with a random number by the client. The variable-length PRG

is instantiated with the stream cipher ChaCha12 [5], which provides 256-bit security. sPIR is instantiated with SealPIR using SEAL v3.2.0 (<https://github.com/ndokmai/sealpir-rust>). We modify the SealPIR library to permit extra operations as required by Synchronized APIR. By default, this version of the library sets the degree of ciphertext polynomial to 2,048, and the size of the coefficients to 54 bits, which provides 128-bit security [6]. sPIR instantiation with OnionPIR [11] will be provided in future work.

The open-source implementation of Synchronized APIR is available at <https://github.com/ndokmai/assisted-pir>.

## 6 Analysis

In this section, we provide theorems and proof sketches for Synchronized APIR correctness and privacy. In addition, we provide an analysis for a loose upper bound for the probability of failure of Synchronized APIR in the event that there are hash collisions. We defer full proofs and an analysis with a tighter bound to future work.

### 6.1 Correctness

We will prove the correctness of Synchronized APIR by first assuming that the universal hash functions in the scheme provide no hash collisions. Then, we provide an analysis for the probability of no hash collisions, which implies an upper bound for the probability of failure of Synchronized APIR.

**Theorem 1** (Synchronized APIR Correctness). Suppose an sPIR-correct scheme  $\Pi^{\text{sPIR}}$  and that the universal hash functions  $\text{H}$ ,  $\text{G}$ , and  $\tilde{\text{H}}$  provide no collisions. Following Definition 4.2, the Synchronized APIR scheme  $\Pi^{\text{SynAPIR}}$  is correct for any server set  $\text{ID} = \{0, \dots, n\}$  and  $\text{AID} = \{1, \dots, n\}$ , collusion threshold  $1 \leq t \leq n$ , databases  $\text{DB}_{\text{id}}$  for all  $\text{id} \in \text{ID}$ ,  $l$  Query sessions, and  $k^{\text{sid}} \in \text{Keys}(\text{DB}_0)$  for all  $\text{sid} \in [l]$ . That is, if

- Synchronization Protocol takes inputs  $(\text{DB}_{\text{id}})_{\text{ID}}, t$  and outputs  $(\text{Cat}_{\text{id}})_{\text{ID}}, (\text{Cat}'_{\text{id}})_{\text{AID}}, (M_{\text{id}})_{\text{ID}}, S$
- Setup Protocol takes inputs  $l, (M_{\text{id}})_{\text{AID}}$  and outputs  $\text{spk}, \text{ssk}, (\text{mq}_{\text{id}}^{\text{sid}})_{\text{sid} \in [l], \text{id} \in \text{AID}}, (\text{tok}^{\text{sid}})_{[l]}$
- For all  $\text{sid} \in [l]$ , Query Protocol for session  $\text{sid}$  takes inputs  $(\text{DB}_{\text{id}})_{\text{ID}}, (M_{\text{id}})_{\text{ID}}, S, (\text{mq}_{\text{id}}^{\text{sid}})_{\text{id} \in \text{AID}}, \text{spk}, \text{ssk}, \text{tok}^{\text{sid}}, k^{\text{sid}}$  and outputs  $a^{\text{sid}}$

then for all  $\text{sid} \in [l]$ ,  $\text{DB}_0[k^{\text{sid}}] = a^{\text{sid}}$ .

*Proof sketch.* First, by assuming that the universal hash functions provide no collisions, we can claim that

- $\tilde{H}(k)$  perfectly represents keyword  $k$ , so  $k$  can be efficiently replaced with  $\tilde{H}(k)$ .
- $H(v)$  perfectly represents value  $v$ , so catalog  $\text{Cat}_{\text{id}}$  perfectly represents database  $\text{DB}_{\text{id}}$ . Therefore tagging the catalogs is the same as tagging the databases themselves.
- $G((k, h))$  perfectly represents  $(k, h)$ , so Catalog Intersection Protocol and Keyword Synchronization Protocol are trivially correct.

Next, we consider the properties of the tags  $(M_{\text{id}})_{\text{ID}}, S$ . By construction of TAG (Algorithm 1), for each keyword  $k \in \text{Keys}(\text{DB}_0)$ ,

1. either  $k \in M_0$  or  $k \in S$  but not both
2. if  $k \in M_0$  and  $k \in M_{\text{id}}$  for any  $\text{id} \in \text{AID}$ , then  $\text{DB}_0[k] = \text{DB}_{\text{id}}[k]$
3.  $\bigcup_{k \in M_0} \bigcup_{\text{id} \in C_k} \{(\text{id}, k)\} = \bigcup_{\text{id} \in \text{AID}} \bigcup_{k \in M_{\text{id}}} \{(\text{id}, k)\}$

Property 3. is a consequence of the fact that by construction,  $M_{\text{id}} = \{k \in M_0 \mid \text{id} \in C_k\}$  and  $C_k = \{\text{id} \in \text{AID} \mid k \in M_{\text{id}}\}$ .

If  $k^{\text{sid}} \in S$ , then by construction of QUERY (Algorithm 5),  $\text{sq}^{\text{sid}}$  targets item  $\text{index}(k^{\text{sid}}, S)$  in the filtered database  $(\text{DB}_0[k])_{k \in S}$  in REPLY (Algorithm 6), which is exactly  $\text{DB}_0[k^{\text{sid}}]$ . By sPIR correctness of  $\Pi^{\text{sPIR}}$ , we conclude that  $a^{\text{sid}} = \text{DB}_0[k^{\text{sid}}]$ .

If  $k^{\text{sid}} \in M_0$ , we consider how the mPIR token  $\text{tok}^{\text{sid}}$  and queries  $\text{mq}_{\text{id}}^{\text{sid}}$  are generated. Consider each bit of  $\text{tok}^{\text{sid}}$  by construction in MTOKEN (Algorithm 2): there exists  $k \in M_0$  such that

$$\text{tok}^{\text{sid}}[\text{index}(k, M_0)] = \bigoplus_{\text{id} \in C_k} \text{mq}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})]$$

where  $C_k$  is a set of id's such that  $\text{DB}_0[k] = \text{DB}_{\text{id}}[k]$  by property 2. stated above. This implies

$$\begin{aligned} & \text{tok}^{\text{sid}}[\text{index}(k, M_0)] \cdot \text{DB}_0[k] \\ &= \bigoplus_{\text{id} \in C_k} \text{mq}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \cdot \text{DB}_0[k] \\ &= \bigoplus_{\text{id} \in C_k} \text{mq}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \cdot \text{DB}_{\text{id}}[k] \end{aligned}$$

And therefore, by the equation above and property 3.,

$$\begin{aligned} & \text{tok}^{\text{sid}} \cdot (\text{DB}_0[k])_{k \in M_0} \\ &= \bigoplus_{k \in M_0} \text{tok}^{\text{sid}}[\text{index}(k, M_0)] \cdot \text{DB}_0[k] \\ &= \bigoplus_{k \in M_0} \bigoplus_{\text{id} \in C_k} \text{mq}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \cdot \text{DB}_{\text{id}}[k] \\ &= \bigoplus_{\text{id} \in \text{AID}} \bigoplus_{k \in M_{\text{id}}} \text{mq}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \cdot \text{DB}_{\text{id}}[k] \\ &= \bigoplus_{\text{id} \in \text{AID}} \text{mq}_{\text{id}}^{\text{sid}} \cdot (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}} \end{aligned}$$

Next, because

$$\text{mq}_0^{\text{sid}}[\text{index}(k^{\text{sid}}, M_0)] = \text{tok}^{\text{sid}}[\text{index}(k^{\text{sid}}, M_0)] \oplus 1$$

by construction of MQUERY (Algorithm 4), we have

$$\begin{aligned} & \bigoplus_{\text{id} \in \text{ID}} \text{mq}_{\text{id}}^{\text{sid}} \cdot (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}} \\ &= \bigoplus_{\text{id} \in \text{ID}} \text{mq}_{\text{id}}^{\text{sid}} \cdot (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}} \\ &= \text{tok}^{\text{sid}} \cdot (\text{DB}_0[k])_{k \in M_0} \oplus (1 \cdot \text{DB}_0[k^{\text{sid}}]) \\ & \quad \bigoplus_{\text{id} \in \text{AID}} \text{mq}_{\text{id}}^{\text{sid}} \cdot (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}} \\ &= \text{DB}_0[k^{\text{sid}}] \end{aligned}$$

Thus proven the correctness of mPIR-Only Query Protocol.

To show that the general case in Query Protocol is correct, we observe that when  $k^{\text{sid}} \in M_0$ ,  $\text{sq}^{\text{sid}}$  targets targets item  $|S|$  by construction of QUERY (Algorithm 5). That is,  $\text{sq}^{\text{sid}}$  targets targets  $\text{mr}_0^{\text{sid}}$  by construction of REPLY (Algorithm 6). By sPIR correctness of  $\Pi^{\text{sPIR}}$ , we conclude that in DECODE (Algorithm 7),  $\text{SDECODE}(\text{ssk}, r^{\text{sid}}) = \text{mr}_0^{\text{sid}}$ . And finally,

$$\begin{aligned} a^{\text{sid}} &= \bigoplus_{\text{id} \in \text{ID}} \text{mr}_{\text{id}}^{\text{sid}} \\ &= \bigoplus_{\text{id} \in \text{ID}} \text{mq}_{\text{id}}^{\text{sid}} \cdot (\text{DB}_{\text{id}}[k])_{k \in M_{\text{id}}} \\ &= \text{DB}_0[k^{\text{sid}}] \end{aligned}$$

by correctness of mPIR-Only Query Protocol.  $\square$

Next, we analyze the probability that there is no hash collision when the universal hash functions are used in the context of Synchronized APIR.

**Theorem 2** (Synchronized APIR Hash Non-collision). Let  $H, G, \tilde{H}$  be random universal hash functions from universal families whose hash values are of length  $l_H, l_G, l_{\tilde{H}}$ , respectively. For any  $\text{ID} = \{0, \dots, n\}$  and databases  $\text{DB}_{\text{id}}, \text{id} \in \text{ID}$ , let  $U := \bigcup_{\text{id} \in \text{ID}} \text{DB}_{\text{id}}$ . We define the following non-collision events:

- $E_1(U) := \forall k, k' \in \text{Keys}(U) : k \neq k' \implies \tilde{\text{H}}(k) \neq \tilde{\text{H}}(k')$
- $E_2(U) := \forall v, v' \in \text{Values}(U) : v \neq v' \implies \text{H}(v) \neq \text{H}(v')$
- $E_3(U) := \forall (k, v), (k', v') \in U : \text{G}((\tilde{\text{H}}(k), \text{H}(v))) \neq \text{G}((\tilde{\text{H}}(k'), \text{H}(v'))) \implies (\text{H}(k), \text{H}(v)) \neq (\text{H}(k'), \text{H}(v'))$

Let  $p(m, d) := e^{-\frac{m(m-1)}{2d}}$ , then

$$\begin{aligned} & \Pr[E_1(U), E_2(U), E_3(U)] \\ & \approx p(|U|, 2^{l_G}) \cdot p(|\text{Keys}(U)|, 2^{l_H}) \cdot p(|\text{Values}(U)|, 2^{l_H}) \end{aligned}$$

*Proof Sketch.* First we consider  $E_1(U)$  and  $E_2(U)$ , which describe the “birthday” problem where no two individuals share the same birthday. For  $E_1(U)$  the number of birthdays is  $2^{l_H}$  and the number of individuals is  $|\text{Keys}(U)|$ ; likewise for  $E_2(U)$ , the number of birthdays is  $2^{l_H}$  and the number of individuals is  $|\text{Values}(U)|$ . By [12], we have

$$\Pr[E_1(U)] \approx p(|\text{Keys}(U)|, 2^{l_H})$$

$$\Pr[E_2(U)] \approx p(|\text{Values}(U)|, 2^{l_H})$$

Given that  $E_1(U)$  and  $E_2(U)$  occur,  $E_3(U)$  also describes the same birthday problem with  $2^{l_G}$  birthdays and  $|U|$  individuals. Therefore,

$$\Pr[E_3(U) \mid E_1(U), E_2(U)] \approx p(|U|, 2^{l_G})$$

Finally, since  $E_1(U)$  and  $E_2(U)$  are independent events, we have

$$\begin{aligned} & \Pr[E_1(U), E_2(U), E_3(U)] \\ & = \Pr[E_3(U) \mid E_1(U), E_2(U)] \Pr[E_1(U), E_2(U)] \\ & = \Pr[E_3(U) \mid E_1(U), E_2(U)] \Pr[E_1(U)] \Pr[E_2(U)] \\ & \approx p(|U|, 2^{l_G}) \cdot p(|\text{Keys}(U)|, 2^{l_H}) \cdot p(|\text{Values}(U)|, 2^{l_H}) \end{aligned}$$

□

What does Theorem 2 imply for Theorem 1? By definition, we know that if events  $E_1(U), E_2(U), E_3(U)$  take place, then Theorem 1 is true with probability 1. That is,

$$\Pr[\Pi^{\text{SynAPIR}} \text{ is correct} \mid E_1(U), E_2(U), E_3(U)] = 1$$

So,

$$\begin{aligned} & \Pr[\Pi^{\text{SynAPIR}} \text{ is correct}] \\ & \geq \Pr[\Pi^{\text{SynAPIR}} \text{ is correct} \mid E_1(U), E_2(U), E_3(U)] \\ & \quad \cdot \Pr[E_1(U), E_2(U), E_3(U)] \\ & = \Pr[E_1(U), E_2(U), E_3(U)] \end{aligned}$$

And therefore

$$\Pr[\Pi^{\text{SynAPIR}} \text{ is incorrect}] \leq 1 - \Pr[E_1(U), E_2(U), E_3(U)]$$

We note that the probability of incorrectness is not per query but per query keyword; that is, the same keyword either succeeds or fails every time from hash collisions. This implies that the amount of failure does not scale with query traffic loads but with the total number of query keywords made.

The upper bound here can be improved because  $E_1(U), E_2(U), E_3(U)$  are not necessary conditions for the correctness of Synchronized APIR, because there are certain hash collisions that simply do not affect the correctness of the scheme. For example, in the catalogs, it is not an issue if  $v \neq v'$  but  $\text{H}(v) = \text{H}(v')$  as long as pairs  $(k, \text{H}(v))$  and  $(k', \text{H}(v'))$  are in the catalogs and  $k \neq k'$ . In future work, we will provide an analysis with a tighter bound taking into account all of these nuances.

## 6.2 Privacy

To prove the privacy of Synchronized APIR, we first provide the privacy definition of Synchronized APIR which directly corresponds to the privacy definition of APIR in Definition 4.3. Then, we show that mPIR queries are in fact pseudorandom given the variable-length PRG, even if some of the PRG seeds are predetermined. Leveraging this fact and the privacy of sPIR scheme, we show that Synchronized APIR is privacy-preserving.

First, let us provide the privacy definition of Synchronized APIR below. Intuitively, Synchronized APIR is privacy-preserving if the attacker is unable to distinguish between queries generated by keyword  $k_0$  and  $k_1$  during one session, even if the attacker has oracle access to all other query sessions and is provided with some of the PRG seeds used to generate the mPIR queries.

**Definition 6.1** ( $(\lambda, t)$ -Synchronized APIR Privacy). Following Definition 4.3, we define the privacy experiment  $\text{PrivA}_{\mathcal{A}, \Pi^{\text{SynAPIR}}, t}(1^\lambda)$  for the Synchronized APIR scheme given a  $\lambda$ -sPIR privacy-preserving scheme  $\Pi^{\text{sPIR}}$  by Definition A.2 and variable-length PRG  $\text{PRG}(s, \ell) \rightarrow r$  below.

1.  $\mathcal{A}$  chooses the number of Query sessions  $l = \text{poly}(\lambda)$ , well-formed catalog  $\text{Cat}_0$  and digests  $(\text{Dig}_{\text{id}})_{\text{AID}}$ , and outputs  $(l, \text{Cat}_0, (\text{Dig}_{\text{id}})_{\text{AID}})$ .  $\mathcal{A}$  chooses a collusion set  $C \subset \text{ID}$  such that  $|C| \leq t$  and sends  $C$  to the oracle  $\mathcal{O}$ .
2. The following steps are followed to generate public and secret parameters:
  - (a) **Synchronization:** Step ② of Catalog Intersection Protocol in Section 5.2.5 is followed with

inputs  $\text{Cat}_0, (\text{Dig}_{\text{id}})_{\text{AID}}$  to produce  $(\text{Cat}'_{\text{id}})_{\text{AID}}$ . Catalogs are tagged by

$$((M_{\text{id}})_{\text{ID}}, S) \leftarrow \text{TAG}(\text{Cat}_0, (\text{Cat}'_{\text{id}})_{\text{AID}}, t)$$

as per step ② of Synchronization Protocol in Section 5.2.1.

- (b) **Setup:** Setup Protocol in Section 5.2.2 is followed with inputs  $(M_{\text{id}})_{\text{ID}}$  and  $l$  and outputs  $\left( \text{spk}, \text{ssk}, \left( \text{tok}^{\text{sid}} \right)_{[l]} \right)$  (the  $\text{mq}$  output is ignored). The PRG seeds  $(s_{\text{id}})_{\text{AID}}$  generated during this step is also saved.

$(\text{spk}, (M_{\text{id}})_{\text{ID}}, S)$  is given to  $\mathcal{A}$  as public parameters.

3. Initialize session ID  $\text{sid} = 1$ .  $\mathcal{A}$  is given oracle access to QUERY in the following way:

- (a) At  $\text{sid} = 1$ ,  $(s_{\text{id}})_{\text{AID}}$  is given to  $\mathcal{O}$ , and  $\mathcal{O}$  gives  $(s_{\text{id}})_{C \setminus \{0\}}$  to  $\mathcal{A}$  (recall that  $s_{\text{id}}$  is a PRG seed).  
(b)  $\mathcal{A}$  chooses and outputs  $k^{\text{sid}} \in \text{Keys}(\text{Cat}_0)$ .  
(c) A query is generated by

$$(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}}) \leftarrow \text{QUERY} \left( M_0, S, \text{spk}, \text{tok}^{\text{sid}}, k^{\text{sid}} \right)$$

following step ⑦ of Query Protocol in Section 5.2.4.  $(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})$  is given to  $\mathcal{O}$ .

- (d) If  $0 \in C$ ,  $\mathcal{O}$  gives  $(\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})$  to  $\mathcal{A}$ ; otherwise,  $\mathcal{O}$  gives an empty value  $\perp$  to  $\mathcal{A}$ .  
(e)  $\text{sid} \leftarrow \text{sid} + 1$

4. During some session  $\text{sid} = \text{sid}^* \leq l$ ,

- (a)  $\mathcal{A}$  chooses  $k_0^{\text{sid}^*}, k_1^{\text{sid}^*} \in \text{Keys}(\text{Cat}_0)$  and outputs  $(k_0^{\text{sid}^*}, k_1^{\text{sid}^*})$ .  
(b) A uniformly random bit is sampled  $b \leftarrow_{\$} \{0, 1\}$ .  
(c) A query is generated

$$\left( \text{mq}_{0,b}^{\text{sid}^*}, \text{sq}_b^{\text{sid}^*} \right) \leftarrow \text{QUERY} \left( M_0, S, \text{spk}, \text{tok}^{\text{sid}^*}, k_b^{\text{sid}^*} \right)$$

and  $(\text{mq}_{0,b}^{\text{sid}^*}, \text{sq}_b^{\text{sid}^*})$  is given to  $\mathcal{O}$ .

- (d) If  $0 \in C$ ,  $\mathcal{O}$  gives  $(\text{mq}_{0,b}^{\text{sid}^*}, \text{sq}_b^{\text{sid}^*})$  to  $\mathcal{A}$ ; otherwise,  $\mathcal{O}$  gives an empty value  $\perp$  to  $\mathcal{A}$ .  
(e)  $\text{sid} \leftarrow \text{sid}^* + 1$

5.  $\mathcal{A}$  is given more oracle access to QUERY until  $\text{sid} = l$ .

6.  $\mathcal{A}$  outputs  $b^* \in \{0, 1\}$ . The output of the experiment is defined to be 1 if  $b^* = b$  and 0 otherwise.

$\Pi^{\text{SynAPIR}}$  is  $(\lambda, t)$ -APIR privacy-preserving for all PPT adversary  $\mathcal{A}$  if there exists a negligible function  $\text{negl}$  such that

$$\Pr [\text{Priva}_{\mathcal{A}, \Pi^{\text{SynAPIR}}, t}(1^\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda)$$

To show that Synchronized APIR satisfies this definition, we first prove that mPIR queries are pseudorandom, even if some of the PRG seeds are predetermined. (The predetermined seeds are indicated by set  $A$  below.) mPIR query pseudorandomness implies that two queries, generated with different keywords, are computationally indistinguishable from one another.

**Theorem 3** (mPIR Query Pseudorandomness). Suppose a variable-length PRG  $\text{PRG}(s, \ell) \rightarrow r$  where  $|s| = \lambda$  and  $|r| = \ell = \text{poly}(\lambda)$ . For any  $\text{ID} := \{0, \dots, n\}$ ,  $\text{AID} := \{1, \dots, n\}$ ,  $1 \leq t \leq n$ , well-formed catalog  $\text{Cat}_0$  and catalog intersections  $\text{Cat}'_{\text{id}}, \text{id} \in \text{AID}$ ,  $l$  Query sessions,  $\text{sid} \in [l]$ ,  $k^{\text{sid}} \in \text{Keys}(\text{Cat}_0)$ ,  $A \subset \text{AID}$  such that  $|A| < t$ , and  $(s_{\text{id}})_A \in \{0, 1\}^{\lambda \cdot |A|}$ ; define

- $\forall \text{id} \in \text{AID} \setminus A : s_{\text{id}} \leftarrow_{\$} \{0, 1\}^\lambda$
- $((M_{\text{id}})_{\text{AID}}, S) = \text{TAG}(\text{Cat}_0, (\text{Cat}'_{\text{id}})_{\text{id}} \in \text{AID})$
- $\forall \text{id} \in \text{AID} : (\text{mq}_{\text{id}}^{\text{sid}})_{\text{sid} \in [l]} = \text{PRG}(s_{\text{id}}, l \cdot |M_{\text{id}}|)$
- $\forall \text{sid} \in [l] : \text{tok}^{\text{sid}} = \text{MTOKEN}((M_{\text{id}})_{\text{ID}}, (\text{mq}_{\text{id}}^{\text{sid}})_{\text{AID}})$
- $\forall \text{sid} \in [l] : \text{mq}_0^{\text{sid}} = \text{MQQUERY}(M_0, \text{tok}^{\text{sid}}, k^{\text{sid}})$
- $\text{aux} = \left( (M_{\text{id}})_{\text{ID}}, A, (s_{\text{id}})_A, (k^{\text{sid}})_{[l]} \right)$

Then for all PPT distinguisher  $\mathcal{D}$ , there exists a negligible function  $\text{negl}$  such that

$$\left| \Pr [\mathcal{D}(r, \text{aux}) = 1] - \Pr \left[ \mathcal{D} \left( (\text{mq}_0^{\text{sid}})_{[l]}, \text{aux} \right) = 1 \right] \right| \leq \text{negl}(\lambda)$$

where  $r \leftarrow_{\$} \{0, 1\}^{l \cdot |M_0|}$  is sampled uniformly at random.

*Proof sketch.* We will use a hybrid argument to prove the theorem as follows.

**Part 1:** For any  $B \subset \text{AID}$  such that  $A \subseteq B$  and  $|B| = t - 1$ , we will slightly modify the definition of  $\text{mq}_{\text{id}}^{\text{sid}}$  above called  $\tilde{\text{mq}}_{\text{id}}^{\text{sid}}$  where if  $\text{id} \in \text{AID} \setminus B$ , then for all  $\text{sid} \in [l]$ ,  $\tilde{\text{mq}}_{\text{id}}^{\text{sid}} \leftarrow_{\$} \{0, 1\}^{l \cdot |M_{\text{id}}|}$  is sampled uniformly at random; everything else remains unchanged i.e.  $\forall \text{id} \in B, \forall \text{sid} \in [l] : \tilde{\text{mq}}_{\text{id}}^{\text{sid}} = \text{mq}_{\text{id}}^{\text{sid}}$ . This results in the new  $\tilde{\text{mq}}_0^{\text{sid}}, \text{sid} \in [l]$ .

We claim that these two ensembles are perfectly indistinguishable

$$\langle r, \text{aux} \rangle \stackrel{P}{\equiv} \langle (\tilde{\text{mq}}_0^{\text{sid}})_{[l]}, \text{aux} \rangle$$



by observing the construction of  $\text{mq}_0^{\text{sid}}$ . Let

$$\tilde{\text{tok}}^{\text{sid}} = \text{MTOKEN} \left( (M_{\text{id}})_{\text{ID}}, \left( \tilde{\text{mq}}_{\text{id}}^{\text{sid}} \right)_{\text{id} \in \text{AID}} \right)$$

and consider each bit of  $\tilde{\text{tok}}^{\text{sid}}$ : for each  $k \in M_0$ ,

$$\begin{aligned} \tilde{\text{tok}}^{\text{sid}}[\text{index}(k, M_0)] &= \bigoplus_{\text{id} \in C_k} \tilde{\text{mq}}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \\ &= \bigoplus_{\text{id} \in C_k \setminus B} \tilde{\text{mq}}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \\ &\quad \oplus \bigoplus_{\text{id} \in C_k \cap B} \tilde{\text{mq}}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})] \end{aligned}$$

Since  $|C_k| = t$  by construction and  $|B| = t - 1$ , we know that  $C_k \setminus B$  is not empty. Because for all  $\text{id} \in C_k \setminus B$ ,  $\tilde{\text{mq}}_{\text{id}}^{\text{sid}}$  is uniformly random by definition, we conclude that  $\bigoplus_{\text{id} \in C_k \setminus B} \tilde{\text{mq}}_{\text{id}}^{\text{sid}}[\text{index}(k, M_{\text{id}})]$  is uniformly random, and so is  $\tilde{\text{tok}}^{\text{sid}}[\text{index}(k, M_0)]$  and  $\tilde{\text{tok}}^{\text{sid}}$  for all  $\text{sid} \in [l]$ . Thus by construction of  $\text{MQQUERY}$ ,  $\tilde{\text{mq}}_0^{\text{sid}}$  is also uniformly random for all  $\text{sid} \in [l]$ . This proves the perfect indistinguishability.

**Part 2:** Given the variable-length PRG, we claim that these two ensembles are computationally indistinguishable

$$\langle (\tilde{\text{mq}}_0^{\text{sid}})_{[l]}, \text{aux} \rangle \stackrel{c}{\equiv} \langle (\text{mq}_0^{\text{sid}})_{[l]}, \text{aux} \rangle$$

via a reduction proof.

Suppose there exists a PPT distinguisher  $\mathcal{D}$  who can distinguish between the two ensembles, we will construct an adversary  $\mathcal{A}$  using  $\mathcal{D}$  as a subroutine to play the following game:

The challenger  $\mathcal{C}$  is tasking  $\mathcal{A}$  to distinguish between a uniformly random string

$$r_0 \leftarrow_{\S} \{0, 1\}^{\sum_{\text{id} \in \text{AID} \setminus B} l \cdot |M_{\text{id}}|}$$

where  $B$  is defined in Part 1, and a pseudorandom string

$$r_1 \leftarrow (\text{PRG}(s_{\text{id}}, l \cdot |M_{\text{id}}|))_{\text{id} \in \text{AID} \setminus B}$$

where each  $s_{\text{id}} \leftarrow_{\S} \{0, 1\}^\lambda$  is sampled uniformly at random. Upon given  $r_b, b \in \{0, 1\}$ ,  $\mathcal{A}$  defines

$$\left( \left( \text{mq}_{\text{id}}^{\text{sid},*} \right)_{\text{sid} \in [l]} \right)_{\text{id} \in \text{AID} \setminus B} := r_b$$

and  $\text{mq}_{\text{id}}^{\text{sid},*} := \text{mq}_{\text{id}}^{\text{sid}}$  for the rest of  $\text{id} \in B$  and  $\text{sid} \in [l]$ . Next,  $\mathcal{A}$  constructs  $\text{mq}_0^{\text{sid},*}$  out of  $\text{mq}_{\text{id}}^{\text{sid},*}, \text{id} \in \text{AID}$  i.e.

$$\forall \text{sid} \in [l] : \text{tok}^{\text{sid},*} = \text{MTOKEN} \left( (M_{\text{id}})_{\text{ID}}, \left( \text{mq}_{\text{id}}^{\text{sid},*} \right)_{\text{AID}} \right)$$

$$\forall \text{sid} \in [l] : \text{mq}_0^{\text{sid},*} = \text{MQQUERY} \left( M_0, \text{tok}^{\text{sid},*}, k^{\text{sid}} \right)$$

and inputs  $((\text{mq}_0^{\text{sid},*})_{[l]}, \text{aux})$  to  $\mathcal{D}$ . If  $\mathcal{D}$  determines the input is the first ensemble, then  $\mathcal{A}$  outputs 0; otherwise  $\mathcal{A}$  outputs 1. By variable-length PRG assumption, we conclude that the advantage of  $\mathcal{A}$  of guessing the correct string is negligible, and so the two ensembles are indistinguishable.

**Part 3:** By “transitivity” of Part 1 and 2,

$$\begin{aligned} \langle r, \text{aux} \rangle &\stackrel{p}{\equiv} \langle (\tilde{\text{mq}}_0^{\text{sid}})_{[l]}, \text{aux} \rangle \stackrel{c}{\equiv} \langle (\text{mq}_0^{\text{sid}})_{[l]}, \text{aux} \rangle \\ &\implies \langle r, \text{aux} \rangle \stackrel{c}{\equiv} \langle (\text{mq}_0^{\text{sid}})_{[l]}, \text{aux} \rangle \end{aligned}$$

thus proven the theorem.  $\square$

Now that we have proven that mPIR are pseudorandom given the variable-length PRG, we will use a hybrid argument to show that the hybridization between mPIR and sPIR (given that the sPIR scheme is privacy-preserving) in Synchronized APIR results in a privacy-preserving scheme, even if some of the PRG seeds are leaked to the attacker.

**Theorem 4** (Synchronized APIR Privacy). Suppose  $\Pi^{\text{sPIR}}$  is  $\lambda$ -sPIR privacy-preserving according to Definition A.2 and  $\text{PRG}(s, l) \rightarrow r$  is a variable-length PRG, then Synchronized APIR scheme  $\Pi^{\text{SynAPIR}}$  is  $(\lambda, t)$ -APIR privacy-preserving according to Definition 6.1.

*Proof sketch.* We first observe that if  $0 \notin C$  in  $\text{PrivA}_{\mathcal{A}, \Pi^{\text{SynAPIR}}, t}$  i.e. main server  $\mathcal{S}_0$  is not compromised, then  $\mathcal{A}$  does not obtain any information related to  $b$ , which implies no advantage in guessing  $b^*$ ; the scheme is therefore trivially privacy-preserving. For the rest of this proof, we hence only focus on the scenario in which  $\mathcal{A}$  has chosen  $0 \in C$ .

We want to show that the view of  $\mathcal{A}$  in  $\text{PrivA}_{\mathcal{A}, \Pi^{\text{SynAPIR}}, t}$  is computationally indistinguishable between when  $b = 0$  and  $b = 1$ . Formally, define the view of  $\mathcal{A}$  in the experiment as

$$\text{outview}_{\mathcal{A}} := \langle (k^{\text{sid}})_{\text{sid} \neq \text{sid}^*}, (k_0^{\text{sid}^*}, k_1^{\text{sid}^*}), \text{auxout} \rangle$$

for the outputs of  $\mathcal{A}$ , where

$$\text{auxout} := (l, \text{Cat}_0, (\text{Dig}_{\text{id}})_{\text{AID}}, C)$$

and

$$\text{inview}_{\mathcal{A}}(b) :=$$

$$\langle (s_{\text{id}})_{C \setminus \{0\}}, (\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})_{\text{sid} \neq \text{sid}^*}, (\text{mq}_{0,b}^{\text{sid}^*}, \text{sq}_b^{\text{sid}^*}), \text{auxin} \rangle$$

for the inputs of  $\mathcal{A}$ , where

$$\text{auxin} := (\text{spk}, (M_{\text{id}})_{\text{ID}}, S)$$

and finally

$$\text{view}_{\mathcal{A}}(b) := \langle \text{outview}_{\mathcal{A}}, \text{inview}_{\mathcal{A}}(b) \rangle$$

We want to show that

$$\text{view}_{\mathcal{A}}(0) \stackrel{c}{\equiv} \text{view}_{\mathcal{A}}(1)$$

with respect to the security parameter  $\lambda$  through the follows steps:

1. Similarly to  $\text{inview}_{\mathcal{A}}(b)$ , we define

$$\widetilde{\text{inview}}_{\mathcal{A}}(b) := \langle (s_{\text{id}})_{C \setminus \{0\}}, (\text{mq}_0^{\text{sid}}, \text{sq}^{\text{sid}})_{\text{sid} \neq \text{sid}^*}, (\text{mq}_{0,b}^{\text{sid}^*}, \tilde{s}\tilde{q}), \text{auxin} \rangle$$

where  $\tilde{s}\tilde{q} \leftarrow \text{SQUERY}(\text{spk}, 0)$ , and

$$\widetilde{\text{view}}_{\mathcal{A}}(b) := \langle \text{outview}_{\mathcal{A}}, \widetilde{\text{inview}}_{\mathcal{A}}(b) \rangle$$

By  $\lambda$ -sPIR privacy assumption, we claim that

$$\text{view}_{\mathcal{A}}(0) \stackrel{c}{\equiv} \widetilde{\text{view}}_{\mathcal{A}}(0)$$

2. By Theorem 3, we claim that

$$\widetilde{\text{view}}_{\mathcal{A}}(0) \stackrel{c}{\equiv} \widetilde{\text{view}}_{\mathcal{A}}(1)$$

because mPIR queries are pseudorandom.

3. Similarly to step 1, by  $\lambda$ -sPIR privacy assumption we claim that

$$\widetilde{\text{view}}_{\mathcal{A}}(1) \stackrel{c}{\equiv} \text{view}_{\mathcal{A}}(1)$$

4. By “transitivity”, we claim that

$$\text{view}_{\mathcal{A}}(0) \stackrel{c}{\equiv} \widetilde{\text{view}}_{\mathcal{A}}(0) \stackrel{c}{\equiv} \widetilde{\text{view}}_{\mathcal{A}}(1) \stackrel{c}{\equiv} \text{view}_{\mathcal{A}}(1)$$

Thus proven the theorem.  $\square$

## 7 Private DNS Query

In this section, we showcase the application of Synchronized APIR to achieve private DNS query for NS (Name Server) records in DNS cache servers. Indeed, the applicability of Synchronized APIR is not limited to NS records, but NS records make a great example because they are relatively stable with the time-to-live (TTL) of 2 days based on the ICANN’s .com zone file, which we will discuss next.

We assume a hypothetical setting in which there are multiple DNS cache servers who keep separate tables for

different top-level domains (TLD) and types of records. Each server gathers the query statistics per each domain from non-private DNS traffic in the past 24 hours, and updates the cache table with the records of most queried domains. The server then builds a Synchronized APIR database out of this cache table. The cache table is assumed to be much smaller than the total number of records for efficiency.

The Synchronized APIR client chooses a group of Synchronized APIR servers: one as the main server and the rest as assisting servers. The parties then participate in the Synchronized APIR scheme.

For the application demonstrated here, the client must inform the servers she wishes to query the .com NS table. In the table, each domain is assumed to link to the collection of all the associated NS records. This implies that in a successful PIR query, the client will receive *all* the NS records for the queried domain, and decides which one to use.

### 7.1 NS Record Dataset

We analyze the .com Generic Top Level Domain (gTLD) zone file provided by ICANN available per request on <https://czds.icann.org>, accessed on June 16, 2022, to gather statistics about NS records. The data analysis is for the purpose of simulating cache tables for the PIR servers, which will be described in the next section. The .com zone file includes all record types, but we filter it for NS records only. Each NS record contains 1) domain name, 2) time-to-live (TTL), 3) class, 4) type, 5) resource record length, and 6) NS domain name. Each domain name may be linked to multiple NS records, and each record are of variable length, depending on how long the NS domain name is. We summarize the statistics in Table 1.

In this table, we draw particular attention to “max. |records| per domain”. “|records| per domain” here means that if a domain name is linked to record 1, record 2, . . . , record  $n$ , then |records| per domain for this domain is |record 1| + |record 2| + . . . + |record  $n$ |. “max.” indicates the maximum of this value across all the domains in the zone file, which is 759 bytes. This implies that the Synchronized APIR database values need to be at least 759 bytes in length to be able to hold all the NS records linked to each domain name for all the domain names. We thus parameterize the database values to be 1,024 bytes in length as a conservative measure.

### 7.2 Data Simulation Method

Although we have access to NS records in the previous section, what remains unknown is how independent DNS cache servers would behave in the real world. Specifically,

Statistics	Values
#records	382M
#domains	159M
avg. #records per domain	2.4
avg.  domain name	17.7 bytes
avg.  records  per domain	73.1 bytes
max.  records  per domain	759 bytes
TTL	2 days

Table 1: **Statistics for NS records in ICANN’s .com gTLD Zone File (M=10<sup>6</sup>).** “#records” indicates the number of records. “#domains” indicates the number of domains. “|records| per domain” specifies total the size in bytes of all records linked each domain, where each NS record is encoded as TTL|CLASS|TYPE|LEN|NSDOMAIN. The TTL’s are of the same value of 172800 across all the records

we need to know what NS records each server is holding in the cache table at a given time to be able to use this table as a database in Synchronized APIR. To simulate the servers’ cache tables for this purpose, we therefore need to make assumptions about 1) the statistical distribution of domain name popularity, and 2) the frequency of DNS queries over a period of time.

For 1), we assume that the popularity of domain name in DNS queries follows a Zipfian distribution [10,15], with the total number of ranks being the total number of domain names with NS records (159 million domains according to Table 1).

For 2), we assume that a cache server collects DNS query statistics from the normal (i.e. non-private) DNS service over a period of time; after which, the server updates the cache table with the most queried domain names during the period. To assume the number of queries (which follow a Zipfian distribution) over a period of time, we consider two scales of DNS operation: *local* and *regional*. The scale of operation implies the scale of traffic loads. In our setting, local means a US city, whereas regional means the entire US.

To get a sense of what the scales look like in the real world, we obtain DNS query statistics from ICANN Managed Root Server (IMRS) accessible on stats.dns.icann.org on June 18, 2022 to obtain the average queries-per-second (QPS) statistics for NS queries in a 24-hour window. For the local scale, we obtain the QPS within the city of Chicago; for the regional setting, we obtain the QPS within the entire US. The average QPS in a 24-hour window for the local setting is 256.3, which accumulates to 22 million queries in the 24-hour cycle. The average QPS in a 24-hour window for the regional setting is 7032.3, which accumulates to 608 million queries in the 24-hour cycle.

We follow these assumptions and parameters to simu-

late 3 cache servers for Synchronized APIR for local and regional experiment: 1 main and 2 assisting servers with collusion threshold of  $t = 2$ . The cache table size, or  $|\text{DB}|$ , for the local experiment is  $2^{13} \sim 8$  thousand domains, and for the regional experiment  $2^{16} \sim 66$  thousand domains. We base the cache table sizes roughly on Cisco’s *Caching DNS Capacity and Performance Guidelines* [1]. To simulate a cache table, we sample domain name ranks from a Zipfian distribution as many times as the number of queries in 24 hours for each server, where the most  $|\text{DB}|$  popular ranks are kept in the cache table. The actual keywords and values in the cache tables are randomly generated; this does not affect the performance of the scheme since the scheme is agnostic of the actual content of database keywords and values.

**Zipfian parameters.** The next question is what appropriate popularity index  $s$  of the Zipfian distribution to use in the simulation. Wang [15] and Jung *et al.* [10] found the popularity index to be 0.98 and 0.91, respectively, in a local DNS setting at the time of the studies. We surmise that the distribution of domain name popularity is always sensitive to time and geography, and therefore there is no single distribution that is universally true and accurate.

Instead, what we aim to demonstrate in our experiments is Synchronized APIR *in the most optimal conditions*. Since Synchronized APIR’s most expensive computational and communication cost corresponds to the number of sPIR items i.e.  $|S|$ , we want to find a Zipfian popularity index  $s$  that minimizes  $|S|$  to demonstrate the most optimal conditions for both the local and regional setting.

To achieve this, we simulate the cache tables for 3 Synchronized APIR servers at varying  $s$  and cache sizes for the local and regional setting ( $2^{12}, 2^{13}, 2^{14}$  domains for local, and  $2^{15}, 2^{16}, 2^{17}$  domains for regional). The results are shown in Figure 6. Here, the optimality is represented by the percentage of sPIR items in the main cache table, i.e.  $|S|/|\text{DB}_0| \times 100\%$ . The results indices that  $s = 1.0$  is an optimal parameter. (Coincidentally, this is close to the 0.98 and 0.91 in Wang and Jung *et al.*’s study, respectively.) We therefore choose  $s = 1.0$  to evaluate Synchronized APIR in the next section.

### 7.3 Evaluation Settings

We summarize the parameters used to evaluate Synchronized APIR in Table 2 and 3.

To evaluate the scheme, we use the implementation specified in Section 5.3. The computing environment is a desktop computer with an Intel i9-10900 CPU and 64 GB of RAM running Ubuntu 20.04. The client and servers are simulated locally in the same computing environment, where each party occupies one CPU core; the

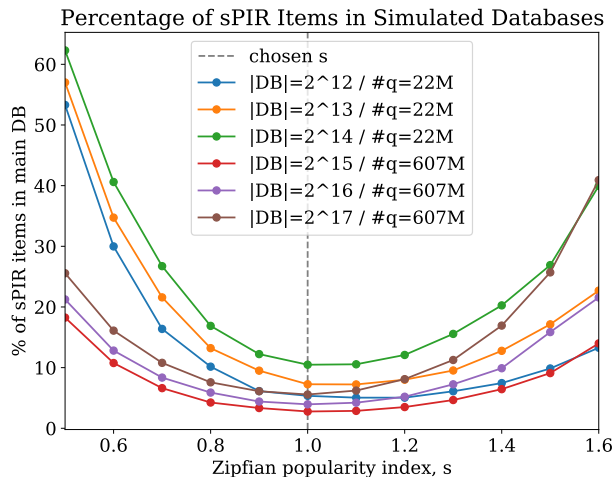


Figure 6: **The effect of Zipfian distribution on the percentage of sPIR items in simulated databases for 3 Synchronized APIR servers at collusion threshold  $t = 2$ .**  $\#q$  denotes the total number of simulated queries.  $s = 1.0$  is chosen as an optimal parameter.

remote communication is via TCP loopback connections.

We measure the performance in terms of communication and computational cost *for a single PIR query*. The communication cost is measured separately in the amount of data uploaded and downloaded by the client in bytes. The computational cost is measured in the amount of CPU time spent by each party in milli-seconds.

Synchronized APIR is evaluated in two experiments: regional and local. In each experiment, a comparison is made between 3 schemes:

1. Baseline SealPIR with query dimensionality  $d = 2$  i.e. recursive PIR
2. Synchronized APIR with  $d = 1$  SealPIR
3. Synchronized APIR with  $d = 2$  SealPIR

The Baseline  $d = 1$  SealPIR scheme is omitted because it repeatedly crashed during the regional experiment in our trial runs regardless of parameters. We suspect this is due to a limitation in the SealPIR library, which is unable to handle databases larger than a certain size for  $d = 1$ .

In Baseline SealPIR experiment, the keyword compression technique is applied to the list of domain names in the cache table; the compressed (i.e. hashed) keywords are sent to the client as the catalog, which allows the client to translate keywords to index positions in the database. In Synchronized APIR, the keyword compression technique is applied per description in Section 5.2.8.

SealPIR is configured by default to provide 128-bit security (see Section 5.3). The log of plaintext modulus is

Parameters	Values
$\#$ servers	3
collusion threshold $t$	2
Zipfian $s$	1.0
$ \text{value} $	1,024 bytes
security parameter	128
SealPIR Parameters	
$\log(\text{plaintext modulus}), d = 1$	14
$\log(\text{plaintext modulus}), d = 2$	16

Table 2: **Shared parameters between local and regional experiments.** “ $\#$ servers” indicates the number of servers.  $|\text{value}|$  indicates the length of each value.  $d$  is the query dimensionality of SealPIR.

Parameters	Local	Regional
$\#$ sampling queries	22M	608M
$ \text{DB} $	$2^{13} \sim 8\text{K}$	$2^{16} \sim 66\text{K}$
$ \text{hash} (\text{bits})$	43	54
Statistics		
% of sPIR	7.3	3.9
$\text{Pr}[\text{failure}]$	$\approx 1.3 \times 10^{-5}$	$\approx 1.2 \times 10^{-5}$

Table 3: **Specific parameters and statistics for local and regional experiments.** “ $\#$ sampling queries” indicates the number of queries to simulate the cache tables. “ $|\text{DB}|$ ” indicates the number of domain names in each cache table. “% of sPIR” indicates the percentage of sPIR items in the main cache table ( $\text{DB}_0$ ). “ $|\text{hash}|$ ” indicates the hash length in bits. “ $\text{Pr}[\text{failure}]$ ” indicates the probability of failure of Synchronized APIR according to Theorem 2.

determines the noise budget in the underlying SEAL fully homomorphic encryption scheme (higher means more noise budget and higher communication cost); this is adjusted by trial and error to ensure the experiments do not fail from the noise after multiple runs.

The sPIR percentage statistics is calculated from the simulated databases. The probability of failure i.e. the probability that there is a hash collision as per Theorem 2 is calculated from the hash length and simulated databases.

## 7.4 Results

The results for the communication costs of the experiments are shown in Table 4 and computational costs in Table 5. We shall refer to the Baseline  $d = 2$  SealPIR scheme as the baseline for comparison.

According to Table 4, the client using Synchronized APIR transmits 4.6x the amount of data that of the base-

Phases	SealPIR, $d = 2$	Sync-APIR, $d = 1$		Sync-APIR, $d = 2$	
	$\mathcal{S}_0$	$\mathcal{S}_0$	$\mathcal{S}_1 + \mathcal{S}_2$	$\mathcal{S}_0$	$\mathcal{S}_1 + \mathcal{S}_2$
<i>Local</i>					
<b>Synchronization</b>					
Download	44K	95K	89K	95K	89K
<b>Setup</b>					
Upload (spk)	3.5M	3.5M	-	3.5M	-
Upload (others)	-	769	19K	769	19K
<b>Query (per session)</b>					
Upload	64K	34K	-	65K	-
Download	257K	32K	2K	257K	2K
<i>Regional</i>					
<b>Synchronization</b>					
Download	352K	799K	704K	799K	704K
<b>Setup</b>					
Upload (spk)	3.5M	3.5M	-	3.5M	-
Upload (others)	-	3.1K	154K	3.1K	154K
<b>Query (per session)</b>					
Upload	64K	40K	-	72K	-
Download	257K	32K	2K	257K	2K

Table 4: **Client’s communication costs for one query in local and regional PIR experiment (in bytes,  $K=2^{10}$ ,  $M=2^{20}$ ).** “ $\mathcal{S}_1 + \mathcal{S}_2$ ” indicates the client’s upload/download costs to/from Server 1 and 2, combined. The upload costs during Setup are listed separately between the cost of SEAL public key and other public parameters.

line in the local experiment and 4.7x in the regional experiment during Synchronization and Setup. The calculation does not include the cost of uploading the long-term SEAL public key, which can be mitigated through public-key infrastructure. During the Query phase of  $d = 1$  and  $d = 2$  Synchronized APIR, the client transmits  $\sim 0.2x$  and  $\sim 1.0x$  the amount of data per query that of the baseline, respectively, in both the local and regional experiment. We note that the client’s upload costs during Query appear irregular with respect to the number of sPIR items due to SealPIR’s query compression technique, which is able to pack many query bits in one large ciphertext; this results in big gaps in query size between small and large number of sPIR query bits.

Where Synchronized APIR clearly outperforms the baseline is in computational cost. According to Table 5, the computational costs of the assisting servers ( $\mathcal{S}_1$  and  $\mathcal{S}_2$ ) are near negligible across the experiments. The computational costs during Synchronization and Setup are slightly cheaper for the main server ( $\mathcal{S}_0$ ) in Synchronized APIR than in the baseline, and slightly more expensive the client in Synchronized APIR than in the baseline. However, the more interesting part is the recurring costs

of the main server during Query. Here, the computational cost of the main server in  $d = 1$  and  $d = 2$  Synchronized APIR are 78% and 12% that of the baseline, respectively.

We summarize our results below.

1. Baseline  $d = 2$  SealPIR offers low one-time communication cost but high recurring costs.
2.  $d = 1$  Synchronized APIR offers 4.6x-4.7x one-time communication cost and  $\sim 0.2x$  recurring communication cost that of the baseline for the client, and 78% recurring computational cost that of the baseline for the main server.
3.  $d = 2$  Synchronized APIR offers 4.6x-4.7x one-time communication cost and  $\sim 1.0x$  recurring communication cost that of the baseline, and 12% recurring computational cost that of the baseline for the main server.

## 7.5 Discussion

It is clear that Synchronized APIR provides significant advantage over SealPIR in a long run after the one-time

Phases	SealPIR, $d = 2$		Sync-APIR, $d = 1$			Sync-APIR, $d = 2$		
	$\mathcal{C}$	$\mathcal{S}_0$	$\mathcal{C}$	$\mathcal{S}_0$	$\mathcal{S}_1$ or $\mathcal{S}_2$	$\mathcal{C}$	$\mathcal{S}_0$	$\mathcal{S}_1$ or $\mathcal{S}_2$
<i>Local</i>								
<b>Synchronization</b>	< 1	2	20	3	3	20	3	3
<b>Setup</b>	107	150	98	89	1	93	86	1
<b>Query (per session)</b>	8	90	1	76	< 1	8	23	< 1
<i>Regional</i>								
<b>Synchronization</b>	7	31	184	35	42	184	35	42
<b>Setup</b>	106	653	221	127	16	214	119	16
<b>Query (per session)</b>	8	414	9	323	4	16	48	4

Table 5: **Computational costs for one query in local and regional PIR experiment (in milli-seconds)**. “ $\mathcal{S}_1$  or  $\mathcal{S}_2$ ” indicates the computational cost on  $\mathcal{S}_1$  or  $\mathcal{S}_2$  since they cost the same CPU time.

cost, but what does this mean in terms of practicality? For a single non-private NS query to a DNS server, the communication cost would be less than 1 KB, and the computational cost negligible, allowing the operation to scale with minimal computational resources. In comparison, Synchronized APIR would be multiple orders of magnitude more expensive. This is not to mention the specificity of the DNS protocol and complexity of the current DNS infrastructure which would require significant adjustment for Synchronized APIR to fit in.

When putting the scheme in context as such, we reckon it would be more sensible to offer private DNS via Synchronized APIR as an optional, special service one needs to opt in to access (especially because our current conception requires the cache servers to collect cache statistics from non-private DNS service). The computational and communication costs of Synchronized APIR are not necessarily prohibitive because 1) the non-sensitive parts of Synchronization and Setup can be outsourced to a third party and shared across multiple clients, allowing the operation to scale, and 2) the recurring cost of a query is only as expensive as the best sPIR scheme available. OnionPIR [11], for example, can reduce the recurring download cost for the client by 25x in comparison to SealPIR at the same computational cost, although a downside is the recurring upload cost may double with a small database.

## 8 Conclusion

In this work, we introduce assisted PIR, a generalization to multi-server PIR that allows for database inconsistencies. We present the construction of Synchronized APIR,

a hybrid APIR protocol between a black-box single-server PIR scheme and a multi-server PIR scheme which takes advantage of the overlap between inconsistent databases to reduce the costs. A formal analysis of Synchronized APIR is also provided.

We apply Synchronized APIR to demonstrate a proof-of-concept private DNS query application, specifically to query NS records, among DNS cache servers. Then, we evaluate the application with simulated datasets based on realistic assumptions about DNS queries and cache behavior.

The results show that despite the higher initial one-time cost, private DNS query via Synchronized APIR is able to outperform the baseline single-server PIR either in communication cost or in computational cost. Although the costs are high in comparison to non-private DNS, Synchronization APIR holds its potential in future developments of single-server PIR schemes from its black-box use of single-server PIR.

## Acknowledgements

We thank Dr. Syed Mahbub Hafiz at University of California, Davis for comments and contributions during the development of this work.

This research was supported in part by the National Science Foundation awards CNS 156537 and the Comcast Innovation Fund.

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the US Government, the National Science Foundation, Comcast, Indiana University, nor University of Calgary.

## References

- [1] Caching DNS capacity and performance guidelines, Nov 2021. Available at [https://www.cisco.com/c/en/us/td/docs/net\\_mgmt/prime/network\\_registrar/10-1/install/guide/Install\\_Guide/Install\\_Guide\\_appendix\\_010001.html](https://www.cisco.com/c/en/us/td/docs/net_mgmt/prime/network_registrar/10-1/install/guide/Install_Guide/Install_Guide_appendix_010001.html).
- [2] Carlos Aguilar-Melchor, Joris Barrier, Laurent Fousse, and Marc-Olivier Killijian. XPIR: Private information retrieval for everyone. *Proceedings on Privacy Enhancing Technologies*, 2(2016):155–174, 2016.
- [3] Andris Ambainis. Upper bound on the communication complexity of private information retrieval. In Pierpaolo Degano, Roberto Gorrieri, and Alberto Marchetti-Spaccamela, editors, *Automata, Languages and Programming*, pages 401–407, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg.
- [4] Sebastian Angel, Hao Chen, Kim Laine, and Srinath Setty. PIR with compressed queries and amortized query processing. In *2018 IEEE symposium on security and privacy (SP)*, pages 962–979. IEEE, 2018.
- [5] Daniel J Bernstein et al. ChaCha, a variant of Salsa20. In *Workshop record of SASC*, volume 8, pages 3–5. Lausanne, Switzerland, 2008.
- [6] Hao Chen, Kyoohyung Han, Zhicong Huang, Amir Jalali, and Kim Laine. Simple encrypted arithmetic library v2.3.0. 2017.
- [7] Benny Chor, Oded Goldreich, Eyal Kushilevitz, and Madhu Sudan. Private information retrieval. In *Proceedings of IEEE 36th Annual Foundations of Computer Science*, pages 41–50. IEEE, 1995.
- [8] David Eppstein, Michael T. Goodrich, Frank Uyeda, and George Varghese. What’s the difference? efficient set reconciliation without prior context. In *Proceedings of the ACM SIGCOMM 2011 Conference*, SIGCOMM ’11, page 218–229, New York, NY, USA, 2011. Association for Computing Machinery.
- [9] Giulia Fanti and Kannan Ramchandran. Efficient private information retrieval over unsynchronized databases. *IEEE Journal of Selected Topics in Signal Processing*, 9(7):1229–1239, 2015.
- [10] Jaeyeon Jung, Emil Sit, Hari Balakrishnan, and Robert Morris. DNS performance and the effectiveness of caching. *IEEE/ACM Transactions on networking*, 10(5):589–603, 2002.
- [11] Muhammad Haris Mughees, Hao Chen, and Ling Ren. OnionPIR: Response efficient single-server pir. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS ’21*, page 2292–2306, New York, NY, USA, 2021. Association for Computing Machinery.
- [12] M. Sayrafiezadeh. The birthday problem revisited. *Mathematics Magazine*, 67(3):220–223, 1994.
- [13] Radu Sion and Bogdan Carbutar. On the computational practicality of private information retrieval. In *Proceedings of the Network and Distributed Systems Security Symposium*, pages 2006–06. Internet Society Geneva, Switzerland, 2007.
- [14] Julien P. Stern. A new and efficient all-or-nothing disclosure of secrets protocol. In Kazuo Ohta and Dingyi Pei, editors, *Advances in Cryptology — ASIACRYPT’98*, pages 357–371, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg.
- [15] Zheng Wang. Analysis of DNS cache effects on query distribution. *The Scientific World Journal*, 2013, 2013.

## A Appendix

### A.1 sPIR Correctness and Privacy Definition

**Definition A.1** (sPIR Correctness). Following Definition 3.1, scheme  $\Pi^{\text{sPIR}}$  is correct for any database  $V = (v_0, \dots, v_{m-1})$  such that  $m \geq 1$  and  $i \in \{0, \dots, m-1\}$  if

- $(\text{spk}, \text{ssk}) \leftarrow \text{SGEN}(1^\lambda)$
- $\text{sq} \leftarrow \text{SQUERY}(\text{spk}, i, m)$
- $\text{sr} \leftarrow \text{SREPLY}(\text{spk}, V, \text{sq})$
- $\text{sa} \leftarrow \text{SDECODE}(\text{ssk}, \text{sr})$

then  $\text{sa} = v_i$ .

**Definition A.2** ( $\lambda$ -sPIR Privacy). Define an sPIR privacy experiment  $\text{PrivS}_{\mathcal{A}, \Pi^{\text{sPIR}}}(1^\lambda)$  for sPIR scheme  $\Pi^{\text{sPIR}}$  according to Definition 3.1 and adversary  $\mathcal{A}$  below.

1.  $\mathcal{A}$  chooses and outputs  $m$ .
2. The parameters are generated  $(\text{spk}, \text{ssk}) \leftarrow \text{SGEN}(1^\lambda)$  and  $\text{spk}$  is given to  $\mathcal{A}$ .
3.  $\mathcal{A}$  is given oracle access to  $\text{SQUERY}(\text{spk}, \cdot)$ .
4.  $\mathcal{A}$  chooses  $i_0^*, i_1^* \in \{0, \dots, m-1\}$  and outputs  $(i_0^*, i_1^*)$ . A uniformly random bit is chosen  $b \leftarrow_{\$} \{0, 1\}$ . A query  $\text{sq}^* \leftarrow \text{SQUERY}(\text{spk}, i_b^*)$  is generated and  $\text{sq}^*$  is given to  $\mathcal{A}$ .

5.  $\mathcal{A}$  is given more oracle access to  $\text{SQUERY}(\text{spk}, \cdot)$ .
6.  $\mathcal{A}$  outputs  $b^* \in \{0, 1\}$ . The output of the experiment is defined to be 1 if  $b^* = b$  and 0 otherwise.

$\Pi^{\text{sPIR}}$  is  $\lambda$ -sPIR privacy-preserving if for all PPT adversary  $\mathcal{A}$ , there exists a negligible function  $\text{negl}$  such that

$$\Pr [\text{PrivS}_{\mathcal{A}, \Pi^{\text{sPIR}}}(1^\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda)$$