

# Communication–Computation Trade-offs in PIR

Asra Ali  
Google LLC  
asraa@google.com

Tancredi Lepoint  
Google LLC  
tancrede@google.com

Sarvar Patel  
Google LLC  
sarvar@google.com

Mariana Raykova  
Google LLC  
mpr2111@columbia.edu

Phillipp Schoppmann  
Humboldt-Universität zu Berlin  
schoppmann@informatik.hu-berlin.  
de

Karn Seth  
Google LLC  
karn@google.com

Kevin Yeo  
Google LLC  
kwlyeo@google.com

## ABSTRACT

In this work, we study the computation and communication costs and their possible trade-offs in various constructions for private information retrieval (PIR), including schemes based on homomorphic encryption (HE) and the Gentry–Ramzan PIR (ICALP’05). First, we introduce new packing and compression techniques which extend the construction of SealPIR (S&P’18), and reduce the communication bandwidth by 70% while preserving essentially the same computation cost. We then present MulPIR, a PIR construction based on homomorphic encryption, which leverages multiplicative homomorphism rather than layered additive homomorphism to implement the recursion steps in PIR. This reduces communication even further, at the cost of an increased computational cost for the server. In particular it eliminates the exponential dependence of PIR communication on the recursion depth due to the ciphertext expansion. Therefore, as a side result, we obtain the first implementation of PIR with full recursion. On the other end of the communication–computation spectrum, we take a closer look at Gentry–Ramzan PIR, a scheme with asymptotically optimal communication rate. Here, the bottleneck is the server’s computation, which we manage to reduce significantly. Our optimizations enable a tunable trade-off between communication and computation, which allows us to reduce server computation by as much as 85%, at the cost of an increased query size. We further show how to efficiently construct PIR for sparse databases. Our constructions support batched queries, as well as symmetric PIR. We implement all of our PIR constructions, and compare their communication and computation overheads with respect to each other and previous work for several application scenarios.

## 1 INTRODUCTION

Accessing public databases often brings privacy concerns for the querier as the query may already reveal sensitive information. For example, queries of medical data can reveal sensitive health information, and access patterns of financial data may leak investment strategies. In settings where such privacy leakage has significant risk, clients may shy away from accessing the database. On the flip side, data providers often do not want access to sensitive client queries, as they could later become a liability for them.

Private information retrieval (PIR) is a cryptographic primitive that aims to address the above question by enabling clients to query a database without revealing any information about their queries to the data owner. While the feasibility of this primitive has been resolved for a long time [16], the search for concretely efficient constructions for practical applications has been an active area of research [5, 6, 21, 26, 27, 32, 38, 50, 68]. In this context, there are several parameters and efficiency measures that characterize a PIR setting and determine what solution might be most suitable for a particular scenario. However, a baseline solution that candidate PIR solutions should improve on is the trivial PIR that returns the whole database to the client.

In this work, we take a deep dive into the setting of PIR where data is stored on a single server. This is the relevant PIR model in practical settings where no additional party is available to assist with the data storage and query execution and one does not wish to trust secure hardware. Non-trivial single server PIR constructions are known to require computational assumptions [45], and such solutions bring significant overheads for both the communication and computation costs compared to information theoretic constructions that are possible in the multi-server setting [23]. While theoretical constructions for PIR [45] achieve poly-logarithmic communication, most efficient single server PIR implementations stop short of this goal and implement only variants of the construction with higher asymptotic communication costs [5, 6, 38, 50].

In this paper, we analyze the communication–computation trade-offs that different PIR construction approaches offer and the hurdles towards achieving the optimal asymptotic communication costs in practice. We present a new PIR construction using somewhat-homomorphic encryption, which leverages multiplicative homomorphism as opposed to layers of additive HE, and improves the communication and computation costs of recursion in existing PIR schemes, enabling for the first time measurements with recursion level beyond three. As an alternative to HE-based PIR, we consider the Gentry–Ramzan PIR construction, which achieves optimal communication but has a high computation overhead. Here, we propose a new client-aided model of computation that allows for a tunable trade-off between communication and computation costs. Our constructions support the asymmetric variant of PIR as well as multi-queries using probabilistic batch codes (PBC). We also provide a new construction for sparse databases, a.k.a, keyword PIR [14],

where the number of database entries is much smaller than the query key domain, and the server’s cost depends only on the actual database size as opposed to the key domain size. We implement our new PIR constructions and compare the communication/computation trade-offs they offer against existing constructions instantiated with different HE schemes. We evaluate the PIR schemes using different database shapes motivated by three applications.

## 1.1 Background

*Efficient Constructions of Single Server PIR.* The most efficient (secure) single server PIR constructions implemented in the recent years [5, 6, 21, 26, 27, 32, 38, 50, 68] are based on homomorphic encryption (HE) techniques and achieve sub-linear communication. The baseline PIR solution (with linear communication complexity) has the client send a selection vector proportional to the database size  $n$  encrypted under additive homomorphic encryption, and has the server return a single encrypted entry by performing  $n$  homomorphic multiplications with a constant and  $n$  homomorphic additions. Sub-linear complexity is achieved by using recursion [65]: the database is viewed as a  $d$ -dimensional database, and the query complexity becomes  $O(d \cdot n^{1/d})$ . Now, for the recursion to work with additive homomorphic encryption schemes, the ciphertext after one level of recursion is viewed as a plaintext in the next layer. In particular, if the additive homomorphic encryption scheme has ciphertext expansion  $F$ , the PIR response will include  $F^{d-1}$  ciphertexts (where, e.g.,  $F \geq 6.4$  in lattice-based schemes, as per [5]). This has limited the recursion depth to  $d \leq 3$  in practice [5, 50].

Along this line of work, there are several papers that present implementations with various resource tradeoffs. Aguilar-Melchor et al. [50] present XPIR with small computation costs but quite large communication costs. On the other hand, another line of work [43, 49] obtain much smaller (almost optimal) communication at the cost of significantly larger computation. In a recent work, Angel, Chen, Laine, and Setty [5], present SealPIR that strikes a better balance in the communication–computation cost. SealPIR requires only slightly more computation than XPIR but uses almost 1000 times less communication than XPIR (but does not achieve the almost optimal rate of the works [43, 49]). SealPIR is instantiated with the FV (lattice-based) homomorphic encryption scheme [28]. It builds upon XPIR [50, 65] and adds a clever query compression technique that reduces the query communication complexity from  $O(dn^{1/d})$  to  $O(d \lceil n^{1/d} / N \rceil)$ , where  $N$  is the number of elements that can be packed in one query ciphertext.

Another known PIR construction that achieves logarithmic communication complexity is the construction of Gentry–Ramzan [34], which does not rely on homomorphic encryption. This PIR construction extends the idea from the work of Cachin et al. [11] which proposes to encode the database  $\{D_i\}_{i \in [n]}$  using the Chinese Remainder Theorem (CRT) representation as  $x \in [n]$  s.t.  $x \equiv D_i \pmod{\pi_i}$  for pairwise coprime moduli  $\{\pi_i\}_{i \in [n]}$ . The query for an element at position  $i$  consists of a group  $\mathbb{G}$  and a generator  $g$  of a subgroup of  $\mathbb{G}$  with order  $q\pi_i$ . The server evaluation of the query computes  $h = g^x$  in  $\mathbb{G}$ , which effectively performs a modular reduction in the exponent to select the component  $D_i \pmod{\pi_i}$  masked with the random value  $q$ . The client recovers the value  $D_i$  by computing the discrete logarithm of  $h$  with base  $g^q$ . The work of Cachin

et al. [11] handled only binary data items, and the Gentry–Ramzan construction [34] shows how to handle larger plaintext domains for the database entries and improves the communication rate to constant. While the resulting construction achieves optimal asymptotic communication rate, it has significant computation costs in several places: the generation of prime numbers needed to instantiate different groups  $\mathbb{G}$  at each query, the computation time at the server exponentiating in the query group  $\mathbb{G}$ , and the decoding which requires computing a discrete logarithm. Because of its computational overhead this PIR construction has been rarely considered as a candidate for implementation and practical applications [19, 20, 54].

In recent years, single server PIR has also been studied in slightly different settings. Two works [10, 12] consider *doubly-efficient PIRs* that attempt to obtain schemes with sub-linear computational costs, but require both significant server overhead and new cryptographic assumptions precluding them from practical applications. Another work [56] introduces the notion of *private stateful information retrieval* where clients store some state over multiple queries. Assuming clients perform enough queries, this scheme obtains both smaller communication and computational costs. In contrast, we build PIR schemes suitable for all settings where clients are stateless and our efficiency guarantees will hold regardless of the number of queries performed by the client.

*Specialized PIR Settings.* Multi-query PIR considers the setting where several PIR queries are executed at the same time. Ishai et al. [40] proposed a construction based on batch codes, which achieves asymptotic improvements in the communication and computation amortized cost multi-query PIR but remains impractical. The work on SealPIR [5] presented a construction based on probabilistic batch codes instantiated with Cuckoo hashing in a similar spirit as private set intersection constructions, which amortizes CPU cost while introducing a small probability of failure ( $\approx 2^{-40}$ ).

PIR for sparse databases, also known as *keyword PIR* [14], considers the setting where the database size is much smaller than its index domain. Chor et al. [14] presented a solution that builds a binary search tree over the items in the database and reduces the computation to a logarithmic number PIR queries for the tree levels. Amortized multi-query PSI techniques [13, 17, 59] could also be viewed as solutions in this setting.

Symmetric PIR (SPIR) [36] extends PIR with additional privacy requirement for the database which guarantees that the querier does not learn anything more than the requested item. SPIR is also known as 1-out-of- $n$  oblivious transfer. Naor and Pinkas [51] provided general transformation from PIR to SPIR using oblivious polynomial evaluation, and there have also been direct constructions [46, 48].

## 1.2 Our Contributions

In this paper, we analyze the exact trade-offs between communication and computation in the context of PIR, and we study the best communication complexity that we can achieve in practice. We present a new PIR construction approach that leverages multiplicative homomorphism and enables new communication–computation trade-offs in HE-based PIR that improve the communication costs in existing implementation. We also consider the PIR protocol of Gentry and Ramzan [34], which has opposite cost characteristics with optimal communication and heavy computation, and present new

techniques that drive the communication-computation trade-off for that construction providing the first implementation of this protocol that scales to databases with millions of elements. Together, our protocols provide various *practical* tradeoffs between computation and communication, which we experimentally evaluate using three application scenarios.

*MulPIR: Leveraging Multiplicative Homomorphism.* As we discussed above, the work of Angel et al. [5] proposed compression techniques that enable them to pack selection vectors in the slots of a homomorphic ciphertext and thus they achieve upload communication of  $O(d\lceil n^{1/d}/N \rceil)$  for a database of size  $n$  using HE with modulus size  $N$  and PIR recursion level  $d$ . While this allows to decrease the upload cost by increasing the recursion level  $d$ , the download communication depends exponentially on the recursion level. The reason for this is the PIR selection algorithm used in this work, which uses layers of HE to implement the partial database selection in different recursion layers. This leads to an explosion of the parameter sizes needed for the outermost encryption.

We propose a different approach (MulPIR) that uses both additive and the *multiplicative* homomorphisms of HE to implement the recursive selection by doing one multiplication of encrypted values per recursion step. This reduces the size of the upload and download together from  $O(d\lceil n^{1/d}/N \rceil + F^{d-1})$  in existing approaches, where  $F$  is the number of plaintexts needed to fit a single HE ciphertext, to  $d \cdot \lceil n^{1/d}/N \rceil \cdot c(d)$ , where  $c(d)$  is the size of a ciphertext that supports  $d$  multiplications (hence, a depth of  $\log d$ ).

We further present new compression techniques that allow us to reduce both the upload and download communication cost leveraging packing multiple selection vectors in the same ciphertext and modulus switching. Inspired by optimization techniques in recent proposals of post-quantum secure encryption schemes [4, 9], we apply “bit dropping” techniques, which guarantee that a ciphertext can decrypt correctly even when omitting some of the least significant bits that account for the noise, to further optimize the concrete communication of our scheme. Our new techniques also allow us to achieve the ideal asymptotic communication complexities for PIR in practice and enable for the first time implementation experiments with recursion level beyond three.

*Gentry–Ramzan PIR: New Efficiency Trade-offs.* The Gentry–Ramzan PIR construction [34] achieves optimal communication complexity for several settings but it pays with significant computational cost. Thus, our contributions focus on ways to reduce this computation overhead, which includes new efficient techniques for encoding the server’s database in CRT form needed for the computation in the scheme, new techniques for fast modular exponentiation needed to answer each query, as well as techniques for client-aided PIR that trade-off between communication and computation.

In this PIR protocol, the server database  $\{D_i\}_{i \in [n]}$  needs to be encoded as  $x = D_i \bmod \pi_i$  for  $i \in [n]$ , where  $\pi_i$  are pairwise coprime integers. A naive application of the Chinese Remainder Theorem requires computation at least quadratic in the size of the database. We leverage a divide-and-conquer modular interpolation algorithm [8] that enables us to achieve computation complexity  $\tilde{O}(n \log^2 n)$ . This technique also allows for pre-computation that can be reused for computations that use the same set of moduli  $\pi_i$ .

The main computation cost on the server side is the modular exponentiation, where the server cannot know the prime factorization of the modulus and thus we cannot use techniques that leverage the factorization to speed-up the computation of the exponentiation. Our approach is to compute the exponentiation as a product of precomputed powers of the generator and to use Straus’s algorithm [66] to do this efficiently. This enables a client-aided technique that allows to improve the server’s computation at the price of additional work at the client. In particular since the precomputed powers of the generator are independent of the exponent, they can be computed at the client who knows the order of the group that it is using for the PIR query and thus can compute exponentiation in this group faster by first reducing the exponent modulo the order of the group. This gives a new way to trade-off computation and communication complexity for the protocol. In Section 6, we show evidence that providing several precomputed powers optimizes the server’s work.

We also apply batching techniques leveraging probabilistic batch codes from Cuckoo hashing [52] for a multi-query setting of Gentry–Ramzan PIR, which provide better scalability for broader sets of the database parameters compared to previous batching approaches [39].

*New Construction for Sparse PIR.* We present a new PIR construction for sparse databases, which provides the client with an answer that either contains the corresponding data if the element is present in the database or is empty, otherwise (see Appendix C). Our construction leverages Cuckoo hashing [52] in a new way inspired by ideas from private set intersection [13, 24, 58, 60] and oblivious RAM [37, 55, 57]. In particular, we observe that we can compress the domain of the database from a large sparse domain to a small dense domain using Cuckoo hashing, which in comparison to regular hashing distributes the items in the hash tables guaranteeing that no collisions occur.

*Comparison and Empirical Evaluation of PIR.* We present a comprehensive comparison of the costs of PIR based on homomorphic encryption. This includes detailed concrete efficiency estimates for the ciphertext size and the computation costs for encryption, decryption and homomorphic operations of different HE schemes. We leverage these estimates to profile the efficiency costs of PIR constructions using the corresponding schemes when instantiated with and without recursion. We further present empirical evaluations of implementations of these PIRs with databases of different shapes (numbers of records and entry sizes). Our benchmarks demonstrate that for the majority of the settings constructions based on lattice based HE constructions, which could also offer multiplicative homomorphism, outperform in computation other additive HE schemes. In terms of communication, additive HE solutions have advantage when the dominant communication cost is the download, e.g., in solutions without recursion for small databases with large entries, since these encryption provides best ratio between plaintext and ciphertext.

We evaluate our new PIR construction, MulPIR, that uses somewhat-homomorphic encryption (SHE) and compare it against SealPIR. MulPIR enables a trade-off of computation for communication,

which reduces the communication of SealPIR by 80% while increasing the computation roughly twice. We also provide the first empirical evaluation of PIR with recursive level beyond three (see Appendix D). Surprisingly, we observe that higher recursion level does not necessarily improve communication. This is due to the fact that lattice-based HE encryptions have a complex relationship between parameters sizes, support for homomorphic operations and number of encryption slots. While recursion improves complexity when the database size increases beyond the number of encryption slots in a ciphertext, increasing the database size requires support for more homomorphic operations, which leads to larger parameters and more slots. In our experiments, the Gentry–Ramzan construction always achieves the best communication complexity but comes with a significant computation cost that can be prohibitive in some settings. However, we show that in terms of *monetary* cost, Gentry–Ramzan can outperform all other PIR approaches considered when database elements are small.

Finally, we apply our construction for keyword PIR to a *password checkup* problem, where a client aims to check if their password is contained in a dataset of leaked passwords, without revealing it to the server. Previous approaches to this problem [67] first reveal a  $k$ -anonymous identifier to the server to reduce the number of candidate passwords to compare against to  $k$ , and then apply a variant of Private Set Intersection to compare the current password against the  $k$  candidates. Our implementations of Gentry–Ramzan and MulPIR enable such lookups with communication *sub-linear* in  $k$ , therefore either enabling better anonymity for the same bandwidth, or same anonymity and smaller bandwidth.

## 2 PRELIMINARIES

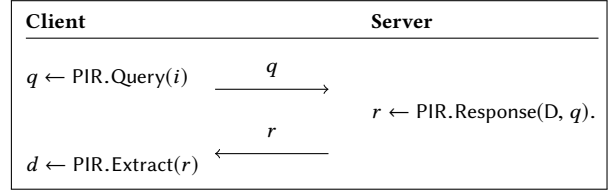
Throughout the rest of this paper, we assume a server owns a database  $D = \{D_1, \dots, D_n\}$  of  $n$  elements, each at most  $l$  bits long.

For any  $m \in \mathbb{Z}$ ,  $m \geq 1$ , we denote by  $[m]$  the interval  $[1, m]$ . We denote by  $\delta_{i,j}$  the Kronecker delta function, defined as  $\delta_{i,j} = 0$  if  $i \neq j$ , and  $\delta_{j,j} = 1$ . For two party computation protocols we will use the notation  $\llbracket a, b \rrbracket$  to denote either inputs or outputs for the two parties, i.e.,  $a$  is either an input or output for the first party, and similarly  $b$  is either input or output for the second party.

### 2.1 Private Information Retrieval (PIR)

*Definition 2.1 (Private Information Retrieval [16]).* A *private information retrieval* protocol addresses the setting where a server holds a database  $D = \{D_1, \dots, D_n\}$  of  $n$  elements, and a client has an input index  $i$ . The goal of the protocol is to enable the client to learn  $D_i$  while guaranteeing that the server does not learn anything about  $i$ . A PIR scheme is specified with the following two algorithms:

- $q \leftarrow \text{PIR.Query}(i)$  – this is an algorithm that the client runs on its input index  $i$  to generate a corresponding query.
- $\llbracket D_i, \perp \rrbracket \leftarrow \text{PIR.Eval}(\llbracket q, D \rrbracket)$  – this is a two-party computation protocol with inputs the client’s encoded query and the server’s database that outputs the corresponding database items to the client. Most PIR constructions are non-interactive and we can replace the evaluation protocol with the following two algorithms (cf. Fig. 1).
  - $r \leftarrow \text{PIR.Response}(D, q)$  – an algorithm that the server runs



**Figure 1: A non-interactive PIR protocol. Correctness of the protocol will ensure that  $d = D_i$ .**

on the client’s encoded query to compute an encoded response.  
–  $D_i \leftarrow \text{PIR.Extract}(r)$  – an algorithm that the client runs on the server’s response to extract the output for the queried item.

*Definition 2.2 (Symmetric Private Information Retrieval (SPIR)).* Symmetric PIR extends the PIR functionality with privacy requirement also for the database guaranteeing the client does not learn anything beyond the element  $D_i$ .

### 2.2 Homomorphic Encryption

For ease of notation and without loss of generality, recall that an additive homomorphic encryption scheme  $\mathcal{HE} = (\text{KeyGen}, \text{Enc}, \text{Dec})$  with plaintext space  $\mathbb{Z}_t$  is an encryption scheme with the following properties:

- $\text{Enc}(\text{sk}, m_1) + \text{Enc}(\text{sk}, m_2) = \text{Enc}(\text{sk}, (m_1 + m_2) \bmod t)$ ,
- $\text{Enc}(\text{sk}, m_1) \cdot \lambda = \text{Enc}(\text{sk}, m_1 \cdot \lambda \bmod t)$ ,

for every  $m_1, m_2, \lambda \in \mathbb{Z}_t$ , for some specific operations  $+$  and  $\cdot$  over the ciphertexts.

Below, we recall the Fan–Vercauteren (FV) homomorphic encryption scheme [28]. For space constraints, the El Gamal and Paillier/Damgård–Jurik cryptosystems are recalled in Appendix A.

*Fan–Vercauteren.* An FV ciphertext is a pair of polynomials over  $R/qR$ , where  $R = \mathbb{Z}[x]/(x^N + 1)$ , and encrypts a message  $m(x) \in R/tR$  for a  $t < q$ . In addition to the standard operations of an encryption scheme (key generation, encryption, decryption), FV also supports homomorphic operations: addition, scalar multiplication, and multiplication.

- **Addition:** Given two ciphertexts  $c_1$  and  $c_2$ , respectively encrypting  $m_1(x)$  and  $m_2(x)$ , the homomorphic addition of  $c_1$  and  $c_2$ , denoted  $c_1 + c_2$ , results in a ciphertext that encrypts the sum  $m_1(x) + m_2(x) \in R/tR$ .
- **Scalar multiplication:** Given a ciphertext  $c \in (R/qR)^2$  encrypting  $m(x) \in R/tR$ , and given  $m'(x) \in R/tR$ , the scalar multiplication of  $c$  by  $m'(x)$ , denoted  $m'(x) \cdot c$ , results in a ciphertext that encrypts  $m'(x) \cdot m(x) \in R/tR$ .
- **Multiplication:** Given two ciphertexts  $c_1$  and  $c_2$ , respectively encrypting  $m_1(x)$  and  $m_2(x)$ , the homomorphic multiplication of  $c_1$  and  $c_2$ , denoted  $c_1 \cdot c_2$ , results in a ciphertext that encrypts the product  $m_1(x) \cdot m_2(x) \in R/tR$ .

Finally, [5] introduced a specific operation called substitution, instantiated using the plaintext slot permutation of [33].

- **Substitution:** Given a ciphertext  $c \in (R/qR)^2$ , that encrypts  $m(x) \in R/tR$ , and an integer  $k$ , the substitution operation

$\text{Sub}_k(\cdot)$  applied on  $c$  results in a ciphertext that encrypts  $m(x^k) \in R/tR$ .

### 3 PIR BASED ON ADDITIVE HE

The majority of private information retrieval constructions that achieve sub-linear communication rely on homomorphic encryption and enable the client to compress its query. More precisely, there are two flavors of homomorphic encryption-based PIR protocols with sub-linear communication that exist in the literature, those based on additive homomorphic encryption (AHE) schemes and those based on fully homomorphic encryption (FHE) schemes.

In this section, we focus on the former flavor, that captures schemes based on El Gamal, Paillier/Damgård–Jurik, and captures the SealPIR protocol proposed by Angel et al. [5] (based on lattice-based additively homomorphic encryption).

We begin by recalling the baseline PIR and the recursion technique in Section 3.1. Then, we recall the SealPIR protocol in Section 3.2. Finally, Section 3.3 presents our contribution and explains how to optimize SealPIR to further reduce the communication by a factor 3 for essentially the same computation cost.

#### 3.1 Background

*Baseline PIR.* We recall the baseline solution for PIR based on homomorphic encryption [45]. Let  $l$  denote the bit-size of the elements of the database and let  $\mathcal{HE} = (\text{KeyGen}, \text{Enc}, \text{Dec})$  be a homomorphic encryption scheme with plaintext space  $\mathbb{Z}_t$  for  $t \geq 2^l$ . Denote by  $C$  the ciphertext space of  $\mathcal{HE}$ . Note that we will interpret each element  $D_i$  as an element of  $\mathbb{Z}_t$ .

The baseline PIR protocol works as follows (cf. Algorithms 1 to 3). To construct the query for index  $k$ , the client encrypts component by component the *selection vector*  $\vec{s} = (s_i)_{i=1\dots n}$  proportional to the size of the database  $n$ , which verifies  $s_i = \delta_{i,k} = 0$  for  $i \neq k$  and  $s_k = \delta_{k,k} = 1$ . To answer the query  $q = (\text{Enc}(\text{sk}, s_i))_{i=1\dots n}$ , the server computes the inner product between the query and the database  $D$  (where  $D_i \in \mathbb{Z}_t$ ), eventually yielding

$$\langle q, D \rangle = \sum_{i=1}^n \text{Enc}(\text{sk}, s_i) \cdot D_i = \text{Enc}\left(\text{sk}, \sum_{i=1}^n \delta_{i,k} D_i\right) = \text{Enc}(\text{sk}, D_k). \quad (1)$$

In the rest of the paper, we will instantiate this protocol with the Paillier/Damgård–Jurik cryptosystem [22, 53], the El-Gamal cryptosystem [31], and an RLWE-based homomorphic cryptosystem [5, 28, 50]. In Appendix A and Table 7, we report on the specific communication and computation costs for the latter schemes.

---

#### Procedure 1 PIR.HE.Query

---

**Input:**  $k \in [1, n]$ .  
 $\vec{s} = (s_i)_{i=1\dots n} = (\delta_{i,k})_{i=1\dots n}$ .  
 $\forall i \in [1, n], q_i \leftarrow \text{Enc}(\text{sk}, s_i)$ .  
**Output:**  $\vec{q} = (q_i)_{i=1\dots n} \in C^n$ .

---



---

#### Procedure 2 PIR.HE.Response

---

**Input:**  $D \in \mathbb{Z}_t^n, \vec{q} \in C^n$ .  
 $r = \langle \vec{q}, D \rangle = \text{Enc}\left(\text{sk}, \langle \vec{s}, D \rangle\right)$  as in Eq. (1).  
**Output:**  $r \in C$ .

---

*Cost of the baseline PIR.* Denote by  $c(n)$  the size of a ciphertext element that enables  $n$  homomorphic scalar multiplications followed by  $n$  homomorphic additions. The overall communication cost is  $n \cdot c(n) + 1 \cdot c(n)$ , hence, is at least linear in the database size. A direct way to trade-off communication of upload and download is by reducing the length  $k$  of the selection vector and returning  $n/k$  database items (assuming  $t > 2^{kl}$ ). This PIR construction could reduce the communication to  $O(n^{1/2})$  ciphertexts, by sending a selection vector of size  $n^{1/2}$  and returning  $O(n^{1/2})$  encrypted database entries.

Two approaches have been proposed in the literature to reduce the overall communication cost: either using recursion (also called folding [32]) using additive homomorphic encryption, or a trivial solution using fully homomorphic encryption. We survey these two approaches below.

*Recursion/Folding.* Kushilevitz, Ostrovsky [45], and later Stern [65], propose the following modification of Algorithms 1 to 3. Instead of representing the database  $D$  as a vector of size  $n$ , one can represent  $D$  as a  $n^{1/2} \times n^{1/2}$  matrix  $M = (M_{i,j})$ , where (say)  $M_{i,j} := D_{in^{1/2}+j}$ . Now, instead of sending (the encryption of) one selection vector  $\vec{s} = (\delta_{i,k})$  of dimension  $n$  for index  $k$ , the client writes  $k = i'n^{1/2} + j'$  where  $i', j' \in [n^{1/2}]$ , and sends two binary selection vectors  $\vec{s}_1 = (s_{1,i}) = (\delta_{i,i'})$  and  $\vec{s}_2 = (s_{2,i}) = (\delta_{j,j'})$  of dimension  $n^{1/2}$ . In particular, it holds that  $s_{1,i} \cdot s_{2,j} = \delta_{i,i'} \cdot \delta_{j,j'} = \delta_{in^{1/2}+j,k}$ , for all  $i, j$ .

The server then performs three steps:

- (1) For each of the  $n^{1/2}$  rows  $M_i = (M_{i,1} \cdots M_{i,n^{1/2}})$ , the server computes the response with the (encryption of the) selection vector  $\vec{s}_2$  as in Eq. (1), i.e., the server obtains the  $n^{1/2}$  ciphertexts

$$c_i = \text{Enc}\left(\text{sk}, \langle \vec{s}_2, (M_{i,j})_j \rangle\right) = \text{Enc}\left(\text{sk}, D_{in^{1/2}+j'}\right).$$

- (2) Since the ciphertext expansion is  $F > 1$ , for each  $i \in [n^{1/2}]$ , the server represents  $c_i$  as  $F$  plaintext elements  $c_{i,1}, \dots, c_{i,F}$ .
- (3) For each of the vectors  $(c_{1,f} \cdots c_{n^{1/2},f})$  with  $f \in [F]$ , the server computes the response with the (encryption of the) selection vector  $\vec{s}_1$  as in Eq. (1), i.e., the server obtains the  $F$  ciphertexts

$$c'_f = \text{Enc}\left(\text{sk}, \langle \vec{s}_1, (c_{i,f})_i \rangle\right) = \text{Enc}\left(\text{sk}, c_{i',f}\right).$$

Upon reception of the response,  $r = (c'_1, \dots, c'_F) \in C^F$ , the client finally extracts the desired result as follows.

- (1) It uses the homomorphic encryption decryption key to recover  $c_{i',f}$  for all  $f \in [F]$ .
- (2) It reconstructs  $c_{i'}$  from the  $c_{i',f}$ 's elements.

---

#### Procedure 3 PIR.HE.Extract

---

**Input:**  $r \in C$ .  
 $d := \text{Dec}(\text{sk}, r) \bmod 2^l$ .  
**Output:**  $d \in \mathbb{Z}_{2^l}$ .

---

<sup>1</sup>We assume without loss of generality that  $F \in \mathbb{Z}$ . Note that we do not ask for any algebraic conditions from the map: for example we could just break down a binary representation of elements of  $C$  into  $F$  plaintexts. For the Paillier cryptosystem, or more precisely the generalization from Damgård and Jurik [22], we will take a different approach: we will select parameters so that the ciphertext after the first folding exactly fits in the plaintext space for the second folding; cf. Appendix A.

---

**Procedure 4** SealPIR.Query

---

**Input:**  $k \in [1, n]$ .

Generate  $\vec{s}_j = (s_{j,i})_{i \in [m]}$  the  $d$  selection vectors in  $\{0, 1\}^m$ .

$\forall j \in [d], m_j \leftarrow \sum_{i \in [m]} s_{j,i} x^i \in R/tR$ .

$\forall j \in [d], q_j \leftarrow \text{Enc}(\text{sk}, m_j)$ .

**Output:**  $\vec{q} = (q_j)_{j \in [d]} \in C^d$ .

---

(3) It uses the homomorphic encryption decryption key on  $c_i$  to recover  $D_{i' n^{1/2+j'}} = D_k$ .

Recursion provides a way to emulate multiplicative homomorphism in one very restricted setting, which, however, suffices for PIR construction. The computation that is enabled by the layering approach for multiplication is inner product with a selection vector that has exactly one non-zero entry which is equal to one. This method easily generalizes by representing the database as a  $d$ -dimensional hyperrectangle  $[n_1] \times \dots \times [n_d]$  with  $n = n_1 \cdot n_2 \cdot \dots \cdot n_d$  (the baseline PIR corresponds to  $d = 1$  with  $n_1 = n$ , and the recursion above to  $d = 2$  with  $n_1 = n_2 = n^{1/2}$ ).

*Cost of recursion.* When  $n_i = n^{1/d}$ , we accomplish the following communication complexity:  $\mathcal{O}(c(n) \cdot dn^{1/d})$  for the user's query and  $\mathcal{O}(F^{d-1}c(n))$  for the server's response. In particular, for small values of  $d$ , we will get sub-linear communication. However, note that for full recursion, i.e.,  $d = \log n$ , communication becomes polynomial in  $n$ .

### 3.2 SealPIR

The SealPIR protocol was proposed by Angel et al. [5], and improves over the XPIR protocol proposed by Aguilar Melchor et al. [50]. Both SealPIR and XPIR instantiate the recursive PIR using the FV homomorphic encryption scheme viewed as an *additive* homomorphic encryption scheme.

In the recursion protocol, the query consists of encryptions of the bits of the selection vectors. The work of [5] starts with the natural observation that many bits can be encrypted in a single ciphertext (and, in particular, at least one per polynomial coefficient) and shows how an encryption with one bit per coefficient can be *obviously expanded* by the server to obtain encryptions of each of the bits in the constant coefficient of the plaintext polynomial. SealPIR's Query algorithm is given in Algorithm 4 and enables to decrease the upload cost by a factor  $\approx N$  (the polynomial ring dimension).<sup>2</sup>

Now, when the server receives such a compressed query, it needs to perform an oblivious expansion into the original query, to then apply Response (Algorithm 2). SealPIR's oblivious expansion is recalled in Algorithm 5.

### 3.3 Optimizing SealPIR

This section explains how to optimize SealPIR. Our techniques yield a reduction of the communication bandwidth by a factor 3x for the same parameters as in [5].

---

<sup>2</sup>Without loss of generality, assume  $m \leq N$ , otherwise the selection vector can be additionally split into  $\lceil m/N \rceil$  different selection vectors. In practice, the selection vectors will be of size  $\leq n^{1/2}$ , which will be below  $N = 2048$  when the database size is  $n \leq 2^{24}$ .

---

**Procedure 5** SealPIR Oblivious Expansion

---

**Input:** Query  $q = \text{Enc}(\sum_{i=0}^{k-1} s_i x^i)$ ,  $k \in [N]$

Find smaller  $m = 2^\ell \geq k$

ciphertexts =  $[q]$

**for**  $j = 0$  to  $\ell - 1$  **do**

**for**  $k = 0$  to  $2^j - 1$  **do**

$c_0 \leftarrow \text{ciphertexts}[k]$

$c_1 \leftarrow x^{-2^j} \cdot c_0$

// scalar multiplication

$c'_k \leftarrow c_0 + \text{Sub}_{N/2^{j+1}}(c_0)$

$c'_{k+2^j} \leftarrow c_1 + \text{Sub}_{N/2^{j+1}}(c_1)$

**end for**

  ciphertexts =  $[c'_0, \dots, c'_{2^{j+1}-1}]$

**end for**

inverse  $\leftarrow m^{-1} \bmod t$

// normalization

**for**  $j = 0$  to  $k - 1$  **do**

$o_j \leftarrow \text{inverse} \cdot \text{ciphertexts}[j]$

**end for**

**Output:** output =  $[o_0, \dots, o_{k-1}]$

---

*Compressing the upload.* Our first contribution comes from the following observation: oblivious expansion (Algorithm 5) is *linear over the plaintext space*. Indeed, all operations used in the algorithms are linear over the plaintext space: additions, substitutions, and scalar multiplications. Hence, it follows that SealPIR's oblivious expansion algorithm enables to expand encryptions of *any* vectors: if  $m = \sum_{i \in [N]} m_i x^i \in R/tR$ , then the output of the oblivious expansion consists of  $N$  ciphertexts, respectively encrypting each of the  $m_i$ 's in the constant coefficient of the plaintexts.

We use the above observation to further compress the size of the query in SealPIR and remark that it could also be applicable in other contexts as well. In SealPIR with recursion  $d$ , the upload consists of  $d \cdot \lceil n^{1/d}/N \rceil$  ciphertexts, where the factor  $d$  comes from the fact that we have  $d$  selection vectors, and  $\lceil n^{1/d}/N \rceil$  comes from the fact that one selection vector of size  $n^{1/d}$  can be embedded in  $\lceil n^{1/d}/N \rceil$  plaintext polynomials in  $R/tR$ . Now we can consider the concatenation of the  $d$  selections vectors of size  $n^{1/d}$  as one vector of size  $d \cdot n^{1/d}$  and use the compression technique over that vector. It follows that the upload size becomes  $\lceil d \cdot n^{1/d}/N \rceil$  ciphertexts. In practice, for  $d \geq 2$ , we usually have  $d \cdot n^{1/d} < N$ , which enables to reduce the upload to a unique ciphertext in SealPIR.

*Compressing the download (and upload).* Our second contribution is to use three compression techniques for homomorphic encryption ciphertexts that will enable to reduce further the communication: secret key encryption, modulus switching, and bit dropping.

- The first optimization comes from the fact that the client, who creates the query ciphertexts, knows the secret key of the homomorphic encryption scheme. In particular, instead of using the public key encryption algorithm, it can use the secret key encryption algorithm of FV. Recall that a FV ciphertext is a tuple  $(c_0, c_1)$  in  $R/qR$ . We briefly describe below the public key and secret key encryption algorithms of FV:

- *Secret Key Encryption.* The secret key is a small polynomial  $s \in R/qR$ . To encrypt  $m \in R/tR$ , sample  $c_0$  uniformly at random in  $R/qR$  and  $e \in R/qR$  a small polynomial, and define  $c_1 = c_0 \cdot s + e + \lceil q/t \rceil m$ .

- *Public Key Encryption.* The secret key is a small polynomial  $s \in R/qR$  and the public key  $(a, b = as + e)$  is an encryption of 0 using the algorithm above. To encrypt  $m \in R/tR$ , sample  $r, e_1, e_2 \in R/qR$  small polynomials, and define  $c_0 = a \cdot r + e_1, c_2 = b \cdot r + e_2 + \lceil q/t \rceil m$ .

A key observation is that when using secret key encryption, the first element  $c_0$  is sampled uniformly at random in  $R/qR$ , whereas it depends on the public key when using public key encryption. Therefore, instead of sending  $c_0$ , the client can instead send a seed  $\rho \in \{0, 1\}^\lambda$ , and the server can reconstruct  $c_0$  from the seed locally. This saves roughly a factor two in size for the upload ciphertexts.

- The second optimization is to use *modulus switching*. This operation allows to transform a ciphertext  $(c_0, c_1) \in (R/qR)^2$  with a noise of norm  $\approx E$  into a ciphertext  $(c_0, c_1) \in (R/pR)^2$  with a noise of norm  $\approx \min(t, (p/q) \cdot E)$  where  $t$  is the plaintext space [18]. It therefore enables us to reduce the download communication in PIR as follows. After finishing to compute the response  $\vec{r} = (r_i)_{i=1 \dots \ell}$  (Algorithm 2), the server will use modulus switching on each ciphertext  $r_i \in (R/qR)^2$  to create a new ciphertext  $r'_i \in (R/pR)^2$ , where  $p \geq t^2$  is chosen large enough to ensure decryption. In practice, this reduces the download size by  $\approx \log_2 q / (2 \log t)$ ; using SealPIR parameters and using modulus switching to a prime  $p \approx 2^{25}$ , this techniques enables to reduce the download by a factor  $60/25 = 2.4$ .
- Finally, we propose to use a technique used in most post-quantum lattice-based encryption schemes proposed for standardization to NIST, such as NewHope [4] and CRYSTAL-Kyber [9], that we call bit-dropping. Essentially, this technique enables to drop the least significant bits of the ciphertext as they carry no information about the message. Indeed, at the end of the PIR computation, each of the ciphertext  $r_i$  in the response if a tuple  $(c_0, c_1) \in (R/qR)^2$  such that

$$c_1 - c_0 \cdot s \bmod q = \lceil q/t \rceil \cdot m + e \in \mathbb{Z},$$

where  $\|e\|$  is small (and in particular,  $\|e\|_\infty \leq \lceil q/t \rceil$ ) and  $m$  is the plaintext. Now, assume that instead of  $c_0$  and  $c_1$ , the server sends the  $\log_2 q - b$  most significant bits from  $c'_0, c'_1$  only. This essentially corresponds to defining  $c'_i = c_i + e_i$  where  $e_i$  is a small noise such that the  $b$  least significant bits of  $c'_i$  are 0, and send  $c''_i = c'_i / 2^b$  to the client. Then, the client can reconstruct the  $c'_i$  and compute

$$c'_1 - c'_0 \cdot s \bmod q = \lceil q/t \rceil \cdot m + (e + e_1 - e_0) \bmod q.$$

Now, if  $e + e_1 - e_0$  is small enough, the last equality will hold over  $\mathbb{Z}$  and the client will be able to decrypt the ciphertext and recover  $m$ . This compression technique can be used both for upload and download, and enable saving a few bits per polynomial coefficient.

*Costs and gains.* All the techniques described above can be use concurrently. We report in Table 1 the gains obtained by using these techniques on SealPIR, where the parameters are chosen so as not to affect security. Additionally, the techniques compressing the ciphertexts need only to be performed on the input and output ciphertexts, effectively adding a negligible computational cost to SealPIR.

## 4 USING HOMOMORPHIC MULTIPLICATION

When PIR was introduced, only additively homomorphic encryption schemes were known. Now, fully (resp. somewhat) homomorphic encryption provides an unbounded (resp. bounded) level of multiplicative homomorphism.

In this section, we first recall in Section 4.1 the generic technique that uses fully homomorphic encryption to construct PIR with optimal communication complexity. Our contributions are presented in the following sections. Sections 4.2 to 4.4 present three flavors of PIR that leverages somewhat homomorphic encryption to achieve new computation-communication trade-offs. Finally, Section 4.5 presents MulPIR, a variant of SealPIR that trades off computation for better communication, especially for higher levels of recursion.

In Table 2, we overview the communication-computation trade-offs in homomorphic encryption based PIR protocols.

### 4.1 Fully Homomorphic Encryption Approach

Assume the homomorphic encryption scheme  $\mathcal{HE}$  is fully homomorphic, i.e., (w.l.o.g. for ease of presentation) there exists a Eval procedure that takes as input ciphertexts  $c_i$  for respective messages  $m_i$  and any function description  $f: \mathbb{Z}_t^\kappa \rightarrow \mathbb{Z}_t$ , and outputs a ciphertext of  $f(m_1, \dots, m_\kappa)$ , which we denote

$$\text{Eval}(\{\text{Enc}(\text{sk}, m_i)\}_{i \in [\kappa]}, f) = \text{Enc}(\text{sk}, f(m_1, \dots, m_\kappa)).$$

A possible approach to computing the selection vector for the PIR query using FHE is based on the following observation: the  $i$ -th bit in the PIR query vector is the output of the equality check between the query index  $k$  and  $i$ . Hence, instead of sending the selection vector  $\vec{s}$ , the client can encrypt each bit  $k_j$  of the index  $k$  and send the resulting  $\kappa = \log n$  ciphertexts to the server. The server then homomorphically computes the selection vector and proceeds as in the baseline PIR construction. This construction achieves communication complexity:  $\mathcal{O}(\log n)$  for the user's query and  $\mathcal{O}(1)$  for the server's response (note that the ciphertext size is independent of the database, hence included in the  $\mathcal{O}$  notation.).

### 4.2 SHE-based solution: Equality Circuit

A first approach consists in implementing the protocol described for fully homomorphic encryption schemes that leverage the observation that, since the values  $k$  and  $i$  have at most  $\kappa = \log n$  bits, the arithmetic circuit for computing equality comparison has multiplicative depth  $\log \kappa = \log \log n$ . Indeed, computing the equality comparison bit for two bit values  $b_1$  and  $b_2$  is equivalent to computing  $1 - (b_1 + b_2 - 2b_1b_2)$  over the integers. Note that in our case only one of the bits coming from the query will be encrypted. Thus, bit equality computation will not require any multiplicative homomorphism. The dominant cost is therefore the multiplication of  $\log n$  encrypted bits, which requires  $\log \log n$  multiplicative degree.

Hence, it suffices to use a somewhat homomorphic encryption that supports  $\log \kappa$  nested multiplications. Then the ciphertext size depends on the size of the database and the communication complexity becomes  $\mathcal{O}(c(n) \log n)$  for the user's query and  $\mathcal{O}(c(n))$  for the server's response.

**Table 1: Gain from our compression techniques (Section 3.3), and of MulPIR (Section 4.5), compared to SealPIR.**

Database size	$n = 2^{18}$			$n = 2^{20}$		
	$d = 1$	$d = 2$	$d = 3$	$d = 1$	$d = 2$	$d = 3$
Recursion						
SealPIR upload (kB)	416	64	96	1664	64	96
SealPIR download (kB)	32	256	2048	32	256	2048
Optimized SealPIR upload (kB)	<b>183</b>	14	14	733	14	14
Optimized SealPIR download (kB)	<b>10</b>	82	655	<b>10</b>	82	655
<i>Total communication wrt SealPIR</i>	<b>0.43×</b>	0.30×	0.31×	<b>0.44×</b>	0.30×	0.31×
MulPIR upload (kB)	<b>183</b>	<b>19</b>	59	<b>733</b>	<b>19</b>	59
MulPIR download (kB)	<b>10</b>	<b>21</b>	43	<b>10</b>	<b>21</b>	43
<i>Total communication wrt SealPIR</i>	<b>0.43×</b>	<b>0.13×</b>	<b>0.04×</b>	<b>0.44×</b>	<b>0.13×</b>	<b>0.04×</b>

For SealPIR, we use the same parameters as in [5, Fig. 9]. The plaintext modulus is fixed to  $t = 2^{12} + 1$ . For the optimizations, we use modulus switching to a prime of 25 bits for SealPIR, drop respectively 5 and 8 bits for upload and download ciphertexts. For MulPIR, the parameters depend on the recursion: for  $d = 2$ , we use  $N = 2048$  and  $\log_2(q) = 80$ ,  $\log_2(p) = 48$ ; for  $d = 3$ , we use  $N = 4096$  and  $\log_2(q) = 120$ ,  $\log_2(p) = 50$ .

**Table 2: Communication-Computation Trade-Off of homomorphic encryption based PIR Protocols.**

Recursion	Total Communication in number of ciphertexts		Approximate computation cost Expressed in homomorphic computation unit: A: addition; S: scalar multiplication; M: multiplication		
	$1 \leq d \leq \log n$	$d = \log(n)$	$1 \leq d \leq \frac{\log n}{\log F}$	$\frac{\log n}{\log F} < d \leq \log n$	$d = \log(n)$
Additive HE	$O\left(dn^{\frac{1}{d}} + F^{d-1}\right)$	$O\left(\log n + F^{\log n-1}\right)$	$n(A + S)$	$n^{\frac{1}{d}} F^{d-1}(A + S)$	$F^{\log n-1}(A + S)$
Somewhat HE	$O\left(dn^{\frac{1}{d}}\right)$	$O(\log n)$	$n(A + S) + n^{\frac{d-1}{d}} M$	$n(A + S) + n^{\frac{d-1}{d}} M$	$n(A + S + M)$
Fully HE	–	$O(\log n)$	–	–	$n \log n M + n(A + S)$

This tables aims at giving an insight on the overall trend but does not reflect accurately the costs; e.g., the communication is indicated in number of ciphertexts while the actual size of the ciphertexts may depend on the database size, and similarly the costs of the homomorphic operations differ between each row.

### 4.3 SHE-based solution: Layered Multiplication

A second approach consists in using recursion, but instead of “emulating” multiplications using additions, use the multiplicative homomorphism of SHE. Using the same notation as in Section 3.1, the PIR protocol becomes as follows. The server performs two steps:

- (1) For each of the  $n^{1/2}$  rows  $M_i = (M_{i,1} \cdots M_{i,n^{1/2}})$ , the server computes the response with the (encryption of the) selection vector  $\vec{s}_2$  as in Eq. (1), i.e., the server obtains the  $n^{1/2}$  ciphertexts

$$c_i = \text{Enc}(\text{sk}, \langle \vec{s}_2, (M_{i,j})_j \rangle) = \text{Enc}(\text{sk}, D_{i n^{1/2+j'}}).$$

- (2) The server now computes the response with the (encryption of the) selection vector  $\vec{s}_1$  using homomorphic multiplication, i.e., the server obtains the ciphertext

$$c = \text{Enc}(\text{sk}, \langle \vec{s}_1, \{D_{i n^{1/2+j'}}\}_i \rangle) = \text{Enc}(\text{sk}, D_{i' n^{1/2+j'}}).$$

Upon reception of the response,  $r = c \in \mathcal{C}$ , the client directly uses the HE decryption key to recover  $D_{i' n^{1/2+j'}} = D_k$ .

Here again, this method easily generalizes by representing the database as a  $d$ -dimensional hyperrectangle  $[n_1] \times \cdots \times [n_d]$  with  $n = n_1 \cdot n_2 \cdots n_d$ . When  $n_i = n^{1/d}$ , we accomplish the following communication complexity:  $O(c(n) \cdot dn^{1/d})$  for the user’s query and  $O(c(n))$  for the server’s response.

### 4.4 SHE-based solution: Selection Vector Reconstruction

Note that the approach of Section 4.3 keeps the layered approach of recursion. In particular, performs *sequentially*  $d$  homomorphic multiplications, effectively requiring the somewhat homomorphic encryption scheme to support circuits of multiplicative depth  $d$ . In particular, for full recursion, this means that the SHE scheme needs to support circuits of depth  $\kappa = \log n$ , which increases the size of the ciphertexts compared to the first approach<sup>3</sup>, where the SHE only required to handle depth  $\log \kappa = \log \log n$ .

We propose below a method that trades communication for computation as follows. First, note that

$$D_{i' n^{1/2+j'}} = \langle \vec{s}_1 \otimes \vec{s}_2, \{D_i\}_{i \in [n]} \rangle,$$

where  $\vec{s}_1 \otimes \vec{s}_2$  is the tensor product of  $\vec{s}_1$  and  $\vec{s}_2$ . More generally, if  $\vec{s}_1, \dots, \vec{s}_d$  denote the selection vectors of dimension  $n^{1/d}$ , such that the indices of the 1 element in  $\vec{s}_i$  is  $j_i$ , then

$$D_{\sum_{j=0}^{d-1} j_i \cdot n^{j/d}} = \langle \vec{s}_1 \otimes \cdots \otimes \vec{s}_d, \{D_i\}_{i \in [n]} \rangle.$$

Hence, this hints to a new protocol, where the client sends the  $d \cdot n^{1/d}$  encryptions of the bits  $s_{j_i, i_j}$  for  $j \in [d], i_j \in [n^{1/d}]$ , the

<sup>3</sup>Indeed, the parameters of somewhat homomorphic encryption schemes scales at least linearly in the multiplicative depth (using techniques called modulus switching or relinearization); hence reducing the multiplicative depth exponentially with also reduce the ciphertext size exponentially.



server computes homomorphically

$$\text{Enc}(\text{sk}, s_{1,i_1} \times \cdots \times s_{d,i_d}), \quad \forall i_1, \dots, i_d \in [n^{1/d}],$$

and then computes the inner product with the original database, as in the baseline PIR (cf. Eq. (1)). Now, note that the latter product can be computed using a binary tree of depth  $\log d$ . For full recursion, i.e.,  $d = \log n$ , the dominant cost in this algorithm is the multiplication of  $d = \log n$  encrypted bits, hence requires  $\log d = \log \log n$  multiplicative degree.

#### 4.5 MulPIR: Putting Everything Together

Sections 4.2 to 4.4 and Table 2 illustrate that homomorphic multiplications enable to reduce the communication of recursion, which becomes a bottleneck for large levels of recursion. We note that the FV homomorphic encryption scheme, that is used in SealPIR, is actually a *somewhat* homomorphic encryption scheme, and the parameters can be chosen to handle an a priori bounded number of multiplications.

We therefore propose MulPIR, which combines the layered multiplication approach from Section 4.3 with the optimizations from Sections 3.2 and 3.3. In particular,

- The Query algorithm is the same as in our optimized variant of SealPIR, described in Section 3.3.
- Upon reception of the query, the server obviously expand the query using SealPIR’s oblivious expansion algorithm;
- Then the server runs the layered multiplication algorithm of Section 4.3;
- Next the server compress the response using modulus-switching and bit-dropping Section 3.3;
- Finally, the client extract the database elements as in SealPIR.

MulPIR trades off computation (higher computational costs for the server) for smaller communication (in total communication, and more particularly for the download communication). We report the communication costs in Table 1 and report the computation costs in Section 6. Finally, in Appendix B, we discuss the one-time communication costs for SealPIR and MulPIR associated with sending the Galois Keys required to perform the Substitute algorithm.

### 5 IMPROVING GENTRY–RAMZAN PIR

An alternative to PIR based on homomorphic encryption is the protocol of Gentry and Ramzan [34], which achieves logarithmic communication and a constant communication rate. While it has been implemented in previous work [19, 20, 54], it is usually dismissed due to its computational complexity [3, 19].

In this section, we describe several optimizations to Gentry–Ramzan PIR that allow us to get a practically efficient implementation. Since the main computation bottleneck for large databases is the server computation (cf. Algorithm 6), we focus on optimizing this part of the protocol. We will first revisit the original protocol [34] (Section 5.1). Then, in Section 5.2, we show how to apply existing techniques [8, 63] to speed up the server setup of Gentry–Ramzan PIR. While this is a one-time setup, it is non-trivial to implement with complexity sub-quadratic in the database size. Finally, in Section 5.3, we show how to speed up the response computation with a novel *client-aided* variant of Gentry–Ramzan PIR, using the

fact that the client can perform modular exponentiations more efficiently since he knows the order of the multiplicative group. This results in an interesting communication–computation trade-off, which we explore in Section 6.

#### 5.1 Gentry–Ramzan PIR

The basic PIR protocol of Gentry and Ramzan [34] works by interpreting the server’s database as a number in a Residue Number System (RNS). That is, given  $n$  coprime integers  $\pi_1, \dots, \pi_n$ , with  $\pi_i \geq 2^l$  for all  $i \in [n]$ , we encode  $D$  as an integer  $E$ , such that

$$E \leq \prod_{i=1}^n \pi_i, \quad \text{and} \quad E \equiv D_i \pmod{\pi_i} \text{ for all } i \in [n]. \quad (2)$$

The existence and uniqueness of  $E$  follows from the Chinese Remainder Theorem, which can also be used to compute  $E$  given  $D$  and all  $\pi_i$ . Observe that (2) implies that we can retrieve the element at index  $i$  by reducing  $E$  modulo  $\pi_i$ . The idea of [34] is to have the server perform this reduction *in the exponent* of a multiplicative group, thus hiding  $i$ . We give the description of the PIR protocol in Algorithms 6 to 8, and refer the reader to [34] for the details.

---

##### Procedure 6 PIR.GR.Query

---

**Input:**  $i \in [n]$ , security parameter  $\lambda$ .

$Q_1 := 2q_1 + 1$  s.t.  $Q_1$  and  $q_1$  are prime and  $\log_2(Q_1) \geq \lambda$ .

$Q_2 := 2q_2\pi_i + 1$  s.t.  $Q_2$  and  $q_2$  are prime and  $\log_2(Q_2) \geq \lambda$ .

$m := Q_1Q_2$ .

$g \leftarrow \mathbb{Z}_m$  s.t.  $|\langle g \rangle| = q_1q_2\pi_i$ .

**Output:**  $(m, g) \in \mathbb{Z} \times \mathbb{Z}_m^*$ .

---



---

##### Procedure 7 PIR.GR.Response

---

**Input:**  $D, (m, g) \in \mathbb{Z} \times \mathbb{Z}_m^*$ .

Encode  $D$  as an integer  $E$  as in Eq. (2).

$g' := g^E \pmod{m}$ .

**Output:**  $g' \in \mathbb{Z}_m^*$ .

---



---

##### Procedure 8 PIR.GR.Extract

---

**Input:**  $g' \in \mathbb{Z}_m^*$ .

$h := g^{q_1q_2}$

$h' := g'^{q_1q_2}$

Solve  $h' = h^d$  for  $d$  using Pohlig–Hellman algorithm.

**Output:**  $d \in \mathbb{Z}_{\pi_i}$ .

---

#### 5.2 Fast Modular Interpolation

Before being able to answer queries, the server must encode the database  $D$  according to Eq. (2). Let  $M = \prod_{i=1}^n \pi_i$  be the product of all moduli, and  $M_k = M/\pi_k = \prod_{i=1, i \neq k}^n \pi_i$ . A naive application of the Chinese Remainder Theorem computes  $E$  as follows:

- (1) For each  $k \in [n]$ , use the extended Euclidean algorithm to compute integers  $a_k, b_k$  such that  $a_k M_k + b_k \pi_k = 1$ .
- (2) Compute  $E = \sum_{k=1}^n D_k a_k M_k = \sum_{k=1}^n D_k a_k \left( \prod_{i=1, i \neq k}^n \pi_i \right)$ .

It is clear that a given modulus  $\pi_k$  divides all summands from Step 2 except the  $k$ -th. Then, using the identity from Step 1, we have  $E \equiv D_k a_k M_k \equiv D_k - D_k b_k \pi_k \equiv D_k \pmod{\pi_k}$  for all  $k \in [n]$ . The problem with that solution is that each  $M_k$  has already size  $\Omega(n)$ . While there are quasi-linear variants of integer multiplication [63] and the extended Euclidean algorithm [64], we have to perform each of those at least  $n$  times, and therefore end up with a total running time of  $\Omega(n^2)$ .

To avoid the quadratic complexity, we rely on the modular interpolation algorithm by Borodin and Moenck [8]. Their main observation is that if we divide our set of moduli  $\pi_i$  evenly into two parts, and call the products of those parts  $M_1$  and  $M_2$ , then the first half of the summands in Step 2 above contains  $M_2$  as a factor, while the other half contains  $M_1$ . Thus,  $M_1$  and  $M_2$  can be factored out of the sum, reducing the computation to two smaller sums and two multiplications:

$$E = M_2 \cdot \left( \sum_{k=1}^{\lfloor n/2 \rfloor} d_k a_k \left( \prod_{i=1, i \neq k}^{\lfloor n/2 \rfloor} \pi_i \right) \right) + M_1 \cdot \left( \sum_{k=\lfloor n/2 \rfloor + 1}^n d_k a_k \left( \prod_{i=\lfloor n/2 \rfloor + 1, i \neq k}^n \pi_i \right) \right).$$

Repeating the above transformation recursively leads to a divide-and-conquer algorithm for modular interpolation, which, using the Schönhage-Strassen integer multiplication [63], has a total running time of  $O(n \log^2 n \log \log n)$  [8]. It relies on the fact that the *supermoduli*  $M_1, M_2$  can be pre-computed, as well as the inverses  $a_k$ . This is especially useful, as we can reuse those for multiple interpolations, as long as the set of moduli  $\pi_i$  remains the same. We will make use of this precomputation when applying our implementation of Gentry–Ramzan PIR to databases with large entries (Section 6.3).

### 5.3 Client-Aided Gentry–Ramzan

As we can see in Algorithm 7, to compute the response to a query, the server has to compute a modular exponentiation, where the exponent encodes the entire database as described in the previous section. Prior work [20] has shown that in practice this step is by far the most expensive part in Gentry–Ramzan PIR.

To speed up the response computation, we rely on the well known fact that one can use Euler’s Theorem to perform modular exponentiations of the form  $g^x \pmod m$  by first reducing the exponent modulo  $\varphi(m) = (Q_1 - 1)(Q_2 - 1)$  and computing

$$g^x \pmod m = g^{x \pmod{\varphi(m)}} \pmod m. \quad (3)$$

While we cannot apply this directly to Algorithm 7 because the server does not (and may not) know  $\varphi(m)$ , the client can use Eq. (3) to perform a part of the server’s computation *without knowing*  $E$ , by pre-computing powers of the generator  $g$ .

Concretely, the server rewrites the large exponent  $E$  according to some base  $b \geq 2$ . Without loss of generality, we know that  $E = E_0 + E_1 b + E_2 b^2 + \dots + E_l b^l$ . It follows that  $g^E = g^{E_0} \cdot (g^b)^{E_1} \cdot (g^{b^2})^{E_2} \dots (g^{b^l})^{E_l}$ . Observe that since  $b$  and  $l$  are public, the client can compute the  $l+1$  values  $g, g^b, g^{b^2}, \dots, g^{b^l}$  without knowing the exponent  $E$ . Furthermore, these  $l$  exponentiations may be efficiently computed by the client using the prime factorization of  $m$  as shown

in Eq. (3). Note that revealing the additional powers of  $g$  to the server does not leak any information, as they could be computed by the server as well, just not as fast. Given these  $l+1$  values, the server’s task reduces to the problem of computing the product of multiple parallel exponentiations. To do this efficiently, one can refer to the survey by Bernstein [7]. For our implementation, we choose Straus’s algorithm [66], a description of which can be found in [42, Alg. 14.88]. In our experiments in Section 6.2, we show that in practice the additional workload on the client insignificant compared to the time needed to generate the prime factors of  $m$ .

## 6 EXPERIMENTAL EVALUATION

In this section, we present experimental results that measure the efficiency of different PIR protocols and illustrate some of the possible tradeoffs that they enable. These results can inform decision making of what is the most appropriate PIR instantiation for a particular application.

### 6.1 Experimental Setup

All our experiments are performed on a desktop computer with a Intel(R) Xeon(R) CPU E5-1650 v4 @ 3.60GHz, 64GB of RAM, running Ubuntu. The parameters of the PIR protocols target 112 bits of security. Unless specified otherwise, the parameters are as follows:

- **El Gamal PIR** with the NIST P-224r1 curve. The plaintext size is chosen to be 4 bytes for fast decryption.
- **Damgård–Jurik PIR**: 1160-bit primes. One ciphertext will therefore encrypt about 290s bytes, for  $s \geq 1$  the Damgård–Jurik parameter.
- **MulPIR**: Polynomials of dimension 2048 and a modulus of 60 bits. The plaintext modulus set to  $t = 2^{12} + 1$ . One ciphertext will therefore encrypt  $12 \cdot 2048/8 = 3072$  bytes.
- **Gentry–Ramzan**: 2048-bit modulus and plaintext block size of  $500 \leq 2048/4$  bits, i.e., 62.5 bytes. When specified as “client-aided”, the client sends 15 generators to the server (cf. Section 5.3).

All the implementations are standalone and rely only on OpenSSL for BigInt and elliptic curve operations.

### 6.2 Baseline Computation Costs

We start with an evaluation of the baseline computation cost of the PIR protocols from Sections 3 and 5.1. Note that in this setting, since we do not use recursion (hence no homomorphic multiplication is performed), SealPIR and MulPIR offer essentially the same performance.

In Table 3, we consider a database of 5000 elements of length 288B (such database was used for evaluation in [6]) and evaluate the client and server costs to setup, create a request, respond to this request, and extract the response. We report communication and computation costs when the database is packed (i.e., the database is reshaped so as to maximize the number of elements in the response; as done in SealPIR [5]). For comparison we also report the costs without packing.

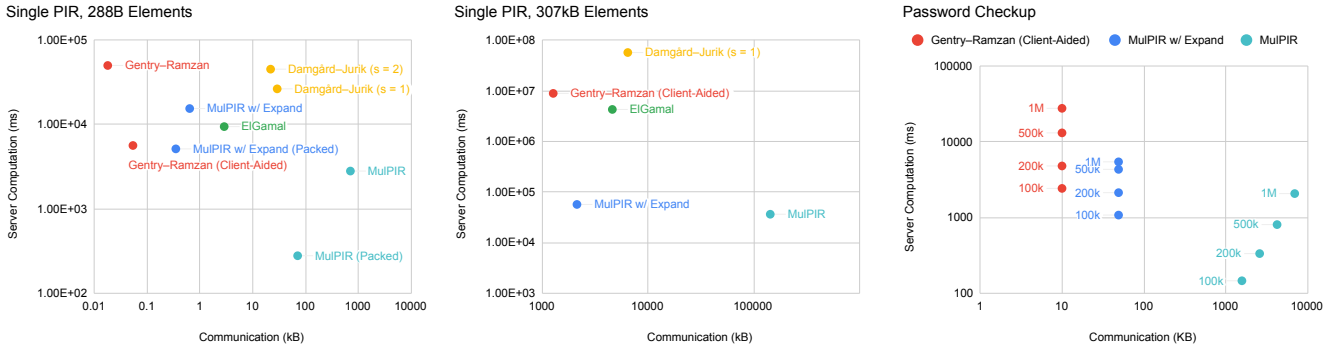
The table also reports on the gain of sending several generators for Gentry–Ramzan. Recall that in Section 5.3, we proposed to use Straus’s algorithm to compute the exponentiation at the core of

**Table 3: Small elements database: 5000 elements of 288B, for a total size of  $\approx$  1MB.**

	packed	Communication (kB)		Computation (ms)					Server Cost (US cents)
		upload	download	C.Setup	S.Setup	C.Create	S.Respond	C.Process	
SealPIR/MulPIR	✗	70480	21	0	522	5474	2803	0.3	0.54
SealPIR/MulPIR with Expand	✗	43	21	0	512	247	15437	0.3	0.0048
SealPIR/MulPIR	✓	7048	21	0	52	550	278	0.3	0.054
SealPIR/MulPIR with Expand	✓	14	21	0	53	242	5136	0.3	0.0017
Damgård–Jurik ( $s = 1$ )	✓	2900	0.6	62806	1	29148	26418	6	0.030
Damgård–Jurik ( $s = 2$ )	✓	2175	0.9	168273	1	32681	45210	14	0.030
Gentry–Ramzan	✓	0.5	1.3	0	1278	12037	49991	361	0.014
Gentry–Ramzan (Client-Aided)	✓	4.1	1.3	0	1280	11327	5631	367	0.0016
ElGamal	✓	280	8	560	22	736	9428	11586	0.0048

Average over 10 computations. “Packed” indicates that the database was reduced to store as many elements as possible per ciphertext. Since Gentry–Ramzan and El Gamal plaintext block sizes are smaller than the size of the entries, respectively 5 and 72 ciphertexts are needed to store a database element. Damgård–Jurik client’s setup includes precomputation to speed up the query creation. Total server costs were computed using Google Cloud Platform prices [1], which were at the time of writing at one cent per CPU-hour and 8 cents per GB of network traffic.

**Figure 2: Trade-off of our PIR implementations, based on data from Table 3 (left), Table 4 (middle), and Table 6 (right).**



the Gentry–Ramzan PIR protocol. Note that the cost of Straus’s algorithm (expressed in number of multiplications in [7] for example), for which one could derive an optimal number of generators to send, does not account for the precomputation cost. However, in Gentry–Ramzan, the server does not know the modulus  $m$  before receiving the client’s request, hence this cost is factored in the server response. In practice, we have determined that about 15 generators was the best communication-computation trade-off one could obtain in Gentry–Ramzan.

Finally, the table shows the price one would have to pay for a single execution of the experiment on Google’s Cloud Platform [1], using a preemptible general-purpose VM with a single CPU core.

### 6.3 Application: Private File Download

Our first application is that of a private file download service. We consider a “fat” database 10,000 files of 307,200 bytes. The total size of the database is therefore 3GB. In this regime, all the PIR protocols are fully packed and need to replicate their operations over “# chunks” ciphertexts. We report communication costs and benchmarks in Table 4.

As expected, Damgård–Jurik and ElGamal are significantly slower than the (packed) variant of SealPIR/MulPIR and Gentry–Ramzan,

respectively for the server and for the client, and will not be considered further in the rest of the section. Additionally, we can see that Expand enables to reduce the communication requirements of SealPIR/MulPIR significantly. While it remains far from the efficient communication cost of Gentry–Ramzan, it offers much better performance.

### 6.4 Application: SealPIR Example Database

This section revisits the application of SealPIR [5], i.e., serving a database of 288B messages. We use this section to further compare the open-source implementation of SealPIR (without modification) available on GitHub [3] against an implementation of MulPIR that uses one homomorphic multiplication (Section 4.5). For this experiment only, and to facilitate comparison, *both* implementations rely on the SEAL homomorphic encryption library [2]—we refer to this MulPIR implementation by “Seal-MulPIR”. We use the same database sizes as in [5] and report the costs in Table 5. The experiment results reflect the use of the more costly homomorphic multiplication in Seal-MulPIR. Note that our custom made implementation of MulPIR (as used in the other sections) will feature a smaller noise growth and hence will enable to select smaller parameters in the Password checkup experiment (Section 6.5).

**Table 4: Private File Download: 10,000 elements of 307kB, for a total database size of  $\approx$  3GB.**

	# chunks	Communication (kB)		Computation (ms)					Total Server Cost (US cents)
		upload	download	C.Setup	S.Setup	C.Create	S.Respond	C.Process	
SealPIR/MulPIR	100	140960	2048	0	105670	11063	36270	25	1.1
SealPIR/MulPIR with Expand	100	71	2048	0	105594	248	56656	25	0.032
Gentry–Ramzan (Client-Aided)	4955	4.1	1259	0	1262133	9676	8848360	344518	2.5
Damgård–Jurik ( $s = 1$ )	1060	5800	614	$\approx$ 60000	$\approx$ 2500	$\approx$ 250000	$\approx$ 57000000	$\approx$ 7000	16
ElGamal	76800	280	4300	$\approx$ 128	$\approx$ 41000	$\approx$ 1500	$\approx$ 4400000	$\approx$ 12500000	1.2

Average over 10 computations. The number of chunks indicates how many ciphertexts are needed to store a database element. The timings indicated with  $\approx$  have been estimated on a smaller number of chunks to finish in a reasonable amount of time. Total server costs were computed as in Table 3.

**Table 5: CPU costs (in ms) of SealPIR and Seal-MulPIR (recursion  $d = 2$ ) for a database of  $n$  elements of 288B.**

	SealPIR ( $d = 2$ ) [3]				Seal-MulPIR ( $d = 2$ )			
	65536	262144	1048576	4194304	65536	262144	1048576	4194304
Database size $n$	65536	262144	1048576	4194304	65536	262144	1048576	4194304
Actual number of rows after packing	6554	26215	104858	419431	6554	26215	104858	419431
Client Setup	40	40	40	40	10	11	11	12
Client Query / Client Extract	1	1	1	1	1	1	1	1
Server Setup	324	1245	4792	19063	768	2982	11772	47819
Server Expand	70	140	279	553	165	330	653	1325
<b>Server Respond</b>	<b>300</b>	<b>907</b>	<b>3087</b>	<b>11513</b>	<b>626</b>	<b>1873</b>	<b>6016</b>	<b>21705</b>

For this comparison only, we re-implemented MulPIR using the Seal library [2] (Seal-MulPIR) and used SealPIR’s implementation from GitHub without modification [3]. Here, because of the noise growth in the Seal library, Seal-MulPIR uses a polynomial dimension of 4096, a 120-bit modulus (product of 2 60-bit moduli), and a plaintext modulus  $t = 2^6 + 1$ . Note that the increased running-time of MulPIR enables the bandwidth savings we reported in Table 1.

## 6.5 Application: Password Checkup

Recent works study the problem of preventing credential stuffing attacks [47, 67] by proposing privacy-preserving protocols where a client queries a centralized breach repository to determine whether her username and password combination has been part of breached data, without revealing the information queried. While this application seems to be a perfect fit for keyword PIR, the size of leaked credentials (4+ billion credentials [67]) remains prohibitively large for PIR. Instead, [47, 67] propose protocols where the client and the server first run an oblivious PRF evaluation (both on usernames and on the tuple username/password), then use the first value to retrieve a bucket and the second value to test for membership after downloading the whole bucket. Precisely, [67] proposes to use  $2^{16}$  buckets, which we infer to contain about  $60k$  elements, and downloading a whole bucket is about 1.6MB of communication.

In this section, we propose to replace the download of the entire bucket with a PIR query. Table 6 shows that using PIR on each bucket is practical (i.e., is comparable to the median waiting time of a few seconds for the client, reported in [67, Tab. 2]) and enables decreasing communication or the number of buckets (or both).

For Gentry–Ramzan, we propose to perform keyword PIR over a bucket using Cuckoo hashing, as introduced in Appendix C.1. The communication is extremely small for any bucket size. For buckets of size 50k, the server computation time is only slightly larger than one second. Unfortunately, the client needs to generate large safe prime numbers which has high computation cost and may impact the applicability of this protocol in practical deployments, such as the one of [67].

Instead, we propose to use MulPIR, which features really low client’s computation costs and low server computation costs. While we could use the Cuckoo hash-based keyword PIR as above, MulPIR would perform worse than Gentry–Ramzan for two reasons. First, the client needs to query as many locations as the number of hash functions. While Gentry–Ramzan supports CRT batching, MulPIR does not support batching natively and its server costs are multiplied by the number of hash functions. Second, a lot of space available in a MulPIR ciphertext is wasted by using Cuckoo hashing, since each bucket row contains at most one element.

Therefore, we propose to use a simpler solution: the server selects a random hash function  $h$  of image size  $k$ , and use it to construct  $k$  bins by placing each of the  $m$  elements  $e$  in the bin of index  $h(e)$ . The client then performs a PIR query over a database of size  $k$ . In order to minimize  $k$ , we want to make the number of elements in each bucket as large as possible while still fitting in one MulPIR ciphertext. Denote  $m = ck \ln k$  for a constant  $c$ . From [62, Th. 1], we know that with overwhelming probability, the maximum size of the bucket will be  $(d_c + 1) \ln k$  where  $d_c$  is the unique root of  $f(x) = 1 + x(\ln c - \ln x + 1) - c$  larger than  $c$ . For every bucket size, we find experimentally the smallest  $k$  such that the whole bin after hashing fits in one MulPIR ciphertext. We report on the communication and computation costs in Table 6. In particular, we conclude that for buckets of size 50k, the server computation time is less than 100ms for about 1MB of communication, and about 1s for about 50kB of communication (plus the one-time keys that need to be transferred), making MulPIR a promising replacement of bucket download in the application of [67].

**Table 6: Password Checkup: Server cost.**

Bucket size	Gentry-Ramzan		MulPIR		MulPIR wo/ Expand	
	Comm. (kB)	Time (ms)	Comm. (kB)	Time (ms)	Comm. (kB)	Time (ms)
10k	<b>10</b>	254	49	540	612	<b>34</b>
20k	<b>10</b>	508	49	540	612	<b>34</b>
50k	<b>10</b>	1308	49	1020	979	<b>69</b>
100k	<b>10</b>	2428	49	1078	1571	<b>146</b>
200k	<b>10</b>	4807	49	2133	2586	<b>334</b>
500k	<b>10</b>	13161	49	4335	4221	<b>807</b>
1M	<b>10</b>	27788	49	5450	6928	<b>2074</b>

The plaintext modulus of MulPIR is  $t = 17$  to enable recursion  $d = 2$ , and  $k$  is respectively equal to 403, 403, 1k, 3k, 8k, 22k, and 58k.

## 7 CONCLUSION

Similar to other advanced cryptographic primitives, PIR is on the verge of transitioning from a theoretical to a practical tool. Our paper presents significant progress in this direction including new PIR constructions and optimization techniques, which provide new ways to trade-off communication and computation. We implement several PIR constructions using different HE schemes as well as the Gentry-Ramzan PIR, and present a comprehensive evaluation of their costs in different settings. Our results demonstrate that the lattice-based FV homomorphic encryption outperforms Paillier and ElGamal in HE-based PIR constructions, while Gentry-Ramzan provides best communication overhead as well as dollar cost for some databases. Our new SHE-based MulPIR enables for the first time an experimental evaluation of full recursion PIR. Overall, our constructions show competitive efficiency for various applications, and we hope our results will serve as a useful reference to inform the choices of PIR construction and parameters in practice.

## REFERENCES

[1] [n.d.]. All Prices | Google Compute Engine Documentation. <https://cloud.google.com/compute/all-pricing>. Accessed 2019-11-01.

[2] 2019. Microsoft SEAL. <https://github.com/microsoft/SEAL>. Accessed 2019-10-30.

[3] 2019. SealPIR: A computational PIR library that achieves low communication costs and high performance. <https://github.com/microsoft/SealPIR>. Accessed 2019-10-30.

[4] Erdem Alkim, Léo Ducas, Thomas Pöppelmann, and Peter Schwabe. 2016. Post-quantum Key Exchange - A New Hope. In *USENIX Security Symposium*. USENIX Association, 327–343.

[5] Sebastian Angel, Hao Chen, Kim Laine, and Srinath T. V. Setty. 2018. PIR with Compressed Queries and Amortized Query Processing. In *IEEE Symposium on Security and Privacy*. IEEE Computer Society, 962–979.

[6] Sebastian Angel and Srinath T. V. Setty. 2016. Unobservable Communication over Fully Untrusted Infrastructure. In *OSDI*. USENIX Association, 551–569.

[7] Daniel J. Bernstein. 2002. Pippenger’s exponentiation algorithm. <http://cr.ypt.to/papers/pippenger.pdf>.

[8] Allan Borodin and R. Moenck. 1974. Fast Modular Transforms. *J. Comput. Syst. Sci.* 8, 3 (1974), 366–386.

[9] Joppe W. Bos, Léo Ducas, Eike Kiltz, Tancrede Lepoint, Vadim Lyubashevsky, John M. Schanck, Peter Schwabe, Gregor Seiler, and Damien Stehlé. 2018. CRYSTALS - Kyber: A CCA-Secure Module-Lattice-Based KEM. In *EuroS&P*. IEEE, 353–367.

[10] Elette Boyle, Yuval Ishai, Rafael Pass, and Mary Wootters. 2017. Can we access a database both locally and privately?. In *Theory of Cryptography Conference*. Springer, 662–693.

[11] Christian Cachin, Silvio Micali, and Markus Stadler. 1999. Computationally Private Information Retrieval with Polylogarithmic Communication. In *Proceedings of the 17th International Conference on Theory and Application of Cryptographic*

*Techniques (EUROCRYPT’99)*.

[12] Ran Canetti, Justin Holmgren, and Silas Richelson. 2017. Towards doubly efficient private information retrieval. In *Theory of Cryptography Conference*. Springer, 694–726.

[13] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. 2018. Labeled PSI from Fully Homomorphic Encryption with Malicious Security. In *ACM Conference on Computer and Communications Security*. ACM, 1223–1237.

[14] Benny Chor, Niv Gilboa, and Moni Naor. 1997. Private Information Retrieval by Keywords.

[15] Benny Chor, Niv Gilboa, and Moni Naor. 1998. Private Information Retrieval by Keywords. *IACR Cryptology ePrint Archive* 1998 (1998), 3.

[16] Benny Chor, Eyal Kushilevitz, Oded Goldreich, and Madhu Sudan. 1998. Private Information Retrieval. *J. ACM* 45, 6 (1998), 965–981.

[17] Michele Ciampi and Claudio Orlandi. 2018. Combining Private Set-Intersection with Secure Two-Party Computation. In *SCN (Lecture Notes in Computer Science)*, Vol. 11035. Springer, 464–482.

[18] Anamaria Costache, Kim Laine, and Rachel Player. 2019. Homomorphic noise growth in practice: comparing BGV and FV. *IACR Cryptology ePrint Archive* 2019 (2019), 493.

[19] Sergiu Costea, Dumitru Marian Barbu, Gabriel Ghinita, and Razvan Rughinis. 2012. A Comparative Evaluation of Private Information Retrieval Techniques in Location-Based Services. In *INCoS*. IEEE, 618–623.

[20] Emiliano De Cristofaro, Yanbin Lu, and Gene Tsudik. 2011. Efficient Techniques for Privacy-Preserving Sharing of Sensitive Information. In *TRUST (Lecture Notes in Computer Science)*, Vol. 6740. Springer, 239–253.

[21] Wei Dai, Yarkin Doröz, and Berk Sunar. 2015. Accelerating SWHE Based PIRs Using GPUs. In *Financial Cryptography Workshops (Lecture Notes in Computer Science)*, Vol. 8976. Springer, 160–171.

[22] Ivan Damgård and Mads Jurik. 2001. A Generalisation, a Simplification and Some Applications of Paillier’s Probabilistic Public-Key System. In *Public Key Cryptography (Lecture Notes in Computer Science)*, Vol. 1992. Springer, 119–136.

[23] Daniel Demmler, Amir Herzberg, and Thomas Schneider. 2014. RAID-PIR: Practical Multi-Server PIR. In *Proceedings of the 6th Edition of the ACM Workshop on Cloud Computing Security (CCSW ’14)*.

[24] Daniel Demmler, Peter Rindal, Mike Rosulek, and Ni Trieu. 2018. PIR-PSI: Scaling private contact discovery. *Proceedings on Privacy Enhancing Technologies* 2018, 4 (2018), 159–178.

[25] Daniel Demmler, Peter Rindal, Mike Rosulek, and Ni Trieu. 2018. PIR-PSI: Scaling Private Contact Discovery. *PoPETs* 2018, 4 (2018), 159–178.

[26] Changyu Dong and Liqun Chen. 2014. A Fast Single Server Private Information Retrieval Protocol with Low Communication Cost. In *ESORICS (1) (Lecture Notes in Computer Science)*, Vol. 8712. Springer, 380–399.

[27] Yarkin Doröz, Berk Sunar, and Ghaith Hammouri. 2014. Bandwidth Efficient PIR from NTRU. In *Financial Cryptography Workshops (Lecture Notes in Computer Science)*, Vol. 8438. Springer, 195–207.

[28] Junfeng Fan and Frederik Vercauteren. 2012. Somewhat Practical Fully Homomorphic Encryption. *IACR Cryptology ePrint Archive* 2012 (2012), 144.

[29] Dimitris Fotakis, Rasmus Pagh, Peter Sanders, and Paul G. Spirakis. 2005. Space Efficient Hash Tables with Worst Case Constant Access Time. *Theory Comput. Syst.* 38, 2 (2005), 229–248.

[30] Michael J. Freedman, Yuval Ishai, Benny Pinkas, and Omer Reingold. 2005. Keyword Search and Oblivious Pseudorandom Functions. In *Proceedings of the Second International Conference on Theory of Cryptography (TCC ’05)*.

[31] Taher El Gamal. 1985. A public key cryptosystem and a signature scheme based on discrete logarithms. *IEEE Trans. Information Theory* 31, 4 (1985), 469–472.

[32] Craig Gentry and Shai Halevi. 2019. Compressible FHE with Applications to PIR. *Cryptology ePrint Archive, Report 2019/733*. <https://eprint.iacr.org/2019/733>.

[33] Craig Gentry, Shai Halevi, and Nigel P. Smart. 2012. Fully Homomorphic Encryption with Polylog Overhead. In *EUROCRYPT (Lecture Notes in Computer Science)*, Vol. 7237. Springer, 465–482.

[34] Craig Gentry and Zulfikar Ramzan. 2005. Single-Database Private Information Retrieval with Constant Communication Rate. In *ICALP (Lecture Notes in Computer Science)*, Vol. 3580. Springer, 803–815.

[35] Yael Gertner, Yuval Ishai, Eyal Kushilevitz, and Tal Malkin. 1998. Protecting Data Privacy in Private Information Retrieval Schemes. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing (STOC ’98)*.

[36] Yael Gertner, Yuval Ishai, Eyal Kushilevitz, and Tal Malkin. 2000. Protecting Data Privacy in Private Information Retrieval Schemes. *J. Comput. Syst. Sci.* 60, 3 (June 2000).

[37] Michael T Goodrich, Michael Mitzenmacher, Olga Ohrimenko, and Roberto Tamassia. 2012. Privacy-preserving group data access via stateless oblivious RAM simulation. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*. Society for Industrial and Applied Mathematics, 157–167.

[38] Matthew Green, Watson Ladd, and Ian Miers. 2016. A Protocol for Privately Reporting Ad Impressions at Scale. In *ACM Conference on Computer and Communications Security*. ACM, 1591–1601.

- [39] Jens Groth, Aggelos Kiayias, and Helger Lipmaa. 2010. Multi-query Computationally-Private Information Retrieval with Constant Communication Rate. In *Public Key Cryptography (Lecture Notes in Computer Science)*, Vol. 6056. Springer, 107–123.
- [40] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. 2004. Batch Codes and Their Applications. In *Proceedings of the Thirty-sixth Annual ACM Symposium on Theory of Computing (STOC '04)*.
- [41] Stanislaw Jarecki and Xiaomin Liu. 2009. Efficient Oblivious Pseudorandom Function with Applications to Adaptive OT and Secure Computation of Set Intersection. In *Proceedings of the 6th Theory of Cryptography Conference on Theory of Cryptography (TCC '09)*.
- [42] Jonathan Katz, Alfred J Menezes, Paul C Van Oorschot, and Scott A Vanstone. 1996. *Handbook of applied cryptography*. CRC press.
- [43] Aggelos Kiayias, Nikos Leonardos, Helger Lipmaa, Kateryna Pavlyk, and Qiang Tang. 2015. Optimal rate private information retrieval from homomorphic encryption. *Proceedings on Privacy Enhancing Technologies* 2015, 2 (2015), 222–243.
- [44] Adam Kirsch, Michael Mitzenmacher, and Udi Wieder. 2009. More Robust Hashing: Cuckoo Hashing with a Stash. *SIAM J. Comput.* 39, 4 (Dec. 2009).
- [45] Eyal Kushilevitz and Rafail Ostrovsky. 1997. Replication is NOT Needed: SINGLE Database, Computationally-Private Information Retrieval. In *FOCS*. IEEE Computer Society, 364–373.
- [46] E. Kushilevitz and R. Ostrovsky. 1997. Replication is Not Needed: Single Database, Computationally-private Information Retrieval. In *Proceedings of the 38th Annual Symposium on Foundations of Computer Science (FOCS '97)*.
- [47] Lucy Li, Bijeeta Pal, Junade Ali, Nick Sullivan, Rahul Chatterjee, and Thomas Ristenpart. 2019. Protocols for Checking Compromised Credentials. In *ACM Conference on Computer and Communications Security*. ACM.
- [48] Helger Lipmaa. 2005. An Oblivious Transfer Protocol with Log-squared Communication. In *Proceedings of the 8th International Conference on Information Security (ISC'05)*.
- [49] Helger Lipmaa and Kateryna Pavlyk. 2017. A simpler rate-optimal CPIR protocol. In *International Conference on Financial Cryptography and Data Security*. Springer, 621–638.
- [50] Carlos Aguilar Melchor, Joris Barrier, Laurent Fousse, and Marc-Olivier Killijian. 2016. XPIR : Private Information Retrieval for Everyone. *PopETs* 2016, 2 (2016), 155–174.
- [51] Moni Naor, Benny Pinkas, and Benny Pinkas. 1999. Oblivious Transfer and Polynomial Evaluation. In *Proceedings of the Thirty-first Annual ACM Symposium on Theory of Computing (STOC '99)*.
- [52] Rasmus Pagh and Flemming Friche Rodler. 2004. Cuckoo Hashing. *J. Algorithms* 51, 2 (May 2004).
- [53] Pascal Paillier. 1999. Public-Key Cryptosystems Based on Composite Degree Residuosity Classes. In *EUROCRYPT (Lecture Notes in Computer Science)*, Vol. 1592. Springer, 223–238.
- [54] Stavros Papadopoulos, Spiridon Bakiras, and Dimitris Papadias. 2012. pCloud: A Distributed System for Practical PIR. *IEEE Trans. Dependable Sec. Comput.* 9, 1 (2012), 115–127.
- [55] Sarvar Patel, Giuseppe Persiano, Mariana Raykova, and Kevin Yeo. 2018. PanORAMA: Oblivious RAM with logarithmic overhead. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE, 871–882.
- [56] Sarvar Patel, Giuseppe Persiano, and Kevin Yeo. 2018. Private stateful information retrieval. In *ACM Conference on Computer and Communications Security*. ACM, 1002–1019.
- [57] Benny Pinkas and Tzachy Reinman. 2010. Oblivious RAM revisited. In *Annual Cryptology Conference*. Springer, 502–519.
- [58] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. 2015. Phasing: Private set intersection using permutation-based hashing. In *24th {USENIX} Security Symposium ({USENIX} Security 15)*. 515–530.
- [59] Benny Pinkas, Thomas Schneider, Oleksandr Tkachenko, and Avishay Yanai. 2019. Efficient Circuit-Based PSI with Linear Communication. In *EUROCRYPT (3) (Lecture Notes in Computer Science)*, Vol. 11478. Springer, 122–153.
- [60] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. 2018. Efficient circuit-based PSI via cuckoo hashing. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 125–157.
- [61] Benny Pinkas, Thomas Schneider, and Michael Zohner. 2018. Scalable Private Set Intersection Based on OT Extension. *ACM Trans. Priv. Secur.* 21, 2 (Jan. 2018).
- [62] Martin Raab and Angelika Steger. 1998. "Balls into Bins" - A Simple and Tight Analysis. In *RANDOM (Lecture Notes in Computer Science)*, Vol. 1518. Springer, 159–170.
- [63] Arnold Schönhage and Volker Strassen. 1971. Schnelle Multiplikation großer Zahlen. *Computing* 7, 3-4 (1971), 281–292.
- [64] Damien Stehlé and Paul Zimmermann. 2004. A Binary Recursive Gcd Algorithm. In *ANTS (Lecture Notes in Computer Science)*, Vol. 3076. Springer, 411–425.
- [65] Julien P. Stern. 1998. A New Efficient All-Or-Nothing Disclosure of Secrets Protocol. In *ASIACRYPT (Lecture Notes in Computer Science)*, Vol. 1514. Springer, 357–371.
- [66] Ernst G. Straus. 1964. Addition chains of vectors (problem 5125). In *American Mathematical Monthly*, Vol. 70. 806–808.
- [67] Kurt Thomas, Jennifer Pullman, Kevin Yeo, Ananth Raghunathan, Patrick Gage Kelley, Luca Invernizzi, Borbala Benko, Tadek Pietraszek, Sarvar Patel, Dan Boneh, and Elie Bursztein. 2019. Protecting accounts from credential stuffing with password breach alerting. In *USENIX Security Symposium*. USENIX Association, 1556–1571.
- [68] Xun Yi, Md. Golam Kaosar, Russell Paulet, and Elisa Bertino. 2013. Single-Database Private Information Retrieval from Fully Homomorphic Encryption. *IEEE Trans. Knowl. Data Eng.* 25, 5 (2013), 1125–1134.

## A APPLICATION TO EXISTING HE SCHEMES

In this section, we discuss instantiations of the PIR approaches from Section 3 with specific homomorphic encryption schemes. In particular, we consider additive ElGamal [31], Paillier/Damgård–Jurik [22], and FV [28], the constructions of which we overview next.

*Additive ElGamal* [31]. Let  $\mathbb{G} = (g)$  be a cyclic group of order  $p$ . The public key is a group element  $h = g^x$ , where the secret key  $x$  is a random integer in  $[p - 1]$ . To encrypt  $m \in [p]$ , sample randomly  $r \leftarrow [p - 1]$  and output  $c = (c_1, c_2) = (g^r, g^m \cdot h^r)$ . To decrypt, compute the discrete logarithm of  $c_2/c_1^x$ . This scheme is additively homomorphic: let  $c = (c_1, c_2)$  encryption  $m$  and  $c' = (c'_1, c'_2)$  encrypting  $m'$ , then  $(c_1 c'_1, c_2 c'_2)$  encrypts  $m_1 + m_2 \bmod p$ . Note that decrypting requires to compute the discrete logarithm in base  $g$ , i.e., we can only decrypt small messages. In particular, an ElGamal ciphertext will have expansion at least  $F \geq 2$ .

*Paillier/Damgård–Jurik* [22]. Let  $N = pq$  be an RSA modulus. The Damgård–Jurik generalization of the Paillier cryptosystem [53] is an additive homomorphic encryption scheme parametrized by an integer  $s$ , such that the plaintext space is  $\mathbb{Z}_{N^s}$  and the ciphertext space is  $\mathbb{Z}_{N^{s+1}}$ . In particular, the ciphertext expansion  $F$  can be made as small as desired and  $F > 1$ . This unusual property enables to simplify the recursion in PIR (cf. Section 3). Using the notation of Section 3, after Step (1), the server obtained  $n^{1/2}$  ciphertexts  $c_i \in \mathbb{Z}_{N^{s+1}}$ . It can then parse this ciphertext as a plaintext element for a Damgård–Jurik scheme with parameter  $s + 1$ ; assuming the selection vector  $\vec{s}_1$  is encrypted under such a scheme, it can then compute  $c' = \text{Enc}_{s+1}(\text{sk}, \langle \vec{s}_1, \{c_i\}_i \rangle) \in \mathbb{Z}_{N^{s+2}}$ . In particular, assume a database with elements in  $N^k$ . The communication after  $d$  levels of recursion, where  $1 \leq d \leq \log n$ , is:

- $n^{1/d}(dk + d(d + 1)/2) \log N$  bits from the client to the server, since each selection vector is encrypted with a modulus  $\log N$  bits larger than the previous one,
- $(d + k) \log N$  bits from the server to the client to send the response.

*FV* [28]. The description of FV is given in Section 2.2; we use the notation of that section. Since FV is additively homomorphic, we can apply the baseline PIR and the recursive PIR protocol of Section 3. The size of a ciphertext is given by  $|ct| = 2N \log q$ . In particular, the communication after  $d$  levels of recursion, for  $1 \leq d \leq \log n$  is

- $(d \cdot n^{1/d} + \lceil 2 \log q / \log t \rceil^d) \cdot (2N \log q)$  bits, from the client to the server where the expansion  $F = 2 \log q / \log t > 2$ ,
- $\lceil 2 \log q / \log t \rceil^{d-1} \cdot (2N \log q)$  bits from the server to the client to send the response.

**Table 7: Bounds on plaintext size, expansion, and decryption cost.**

Scheme	Plaintext size	Expansion $F$	Decryption cost
ElGamal	pt small	$F \geq 2$	$2^{\text{pt}}$ mults of $2^{\lambda_{EG}}$ -bit nums
Damgård–Jurik	$\text{pt} \leq s \cdot \lambda_{DJ}$ bits	$F \geq 1 + 1/s$	1 exponentiation with $\lambda_{DJ}$ -bit exp
FV	$\text{pt} < \log(q) \cdot \lambda_{FV}$ bits	$F \geq 2$	one add and one mult in $\mathbb{Z}_q[x]/(x^{\lambda_{FV}} + 1)$
Gentry–Ramzan	$\text{pt} < \lambda_{GR}/4$	$F > 4$	$4\text{pt} \sqrt{n}$

Here,  $s$  is an integer parameter, and  $\lambda_{DJ}$ ,  $\lambda_{EG}$ ,  $\lambda_{FV}$ , and  $\lambda_{GR}$  are the security parameters for the different encryption schemes, the size of which is determined by the underlying hardness assumptions. Although not exactly an encryption scheme, we include Gentry–Ramzan here. In this case, Decryption corresponds to solving a discrete logarithm, for which the running time depends on the database size  $n$  [34, p. 808].

However, since FV is also somewhat (and fully) homomorphic, we can apply the PIR protocols of Section 4.1. This enables to reduce the communication to

- $(d \cdot n^{1/d}) \cdot (2N \log q)$  bits, from the client to the server,
- $2N \log q$  bits from the server to the client to send the response.

## B (ONE-TIME) KEY INFORMATION SIZE

The communication costs both in Table 1 as well as the SealPIR paper are the communication costs *per query*, assuming the server knows the Galois Keys that will be required to perform the Substitute algorithm that is used in the oblivious expansion algorithm. Similarly for MulPIR, we also require the client to send one additional key-switching key to perform the homomorphic multiplication. Note that all this key information does not depend on the index that is queried, and can be generated beforehand/offline by the client, and reused for multiple query. The communication cost of such key information are provided in Table 9.

Note that SealPIR requires to send  $\log N$  Galois keys, where each Galois key consists of  $\log_2 q/3$  ciphertexts; hence it is possible to use two of the optimizations from Section 3.3: sending a seed rather than a random polynomial, and bit-dropping (in practice  $b = 5$  bits are dropped). Note that the size of the Galois keys in MulPIR may be higher than in SealPIR. Indeed, as explained in Table 1, the number of coefficients and size of moduli depends on the recursion depth.

We propose an optimization that trades computation for communication, as follows. Instead of sending  $\log N$  Galois keys, it is possible to send one Galois keys only (a generator of the Galois group) and apply it repeatedly. For example, for any substitution  $m(x) \mapsto m(x^{2^j+1})$ ,  $j \leq \log N$  that we need to perform during oblivious expansion, the substitution  $m(x) \mapsto m(x^5)$  can be applied repeatedly to get all possible substitution powers. This enables to reduce the number of keys to send from  $\log N$  to 1 Galois key.

## C BEYOND PIR: SPARSITY AND DATABASE PRIVACY

In this section we consider functionalities beyond the traditional setup for PIR that bring extended computation capability, efficiency and security properties, which can be advantageous in different application scenarios.

### C.1 Keyword PIR using Cuckoo Hashing

The traditional setup for PIR constructions assumes that the database entries have public indices which are known to the client submitting queries. In particular, these indices coincide with the domain of all possible queries for the client. Under this assumption the size of the database is equal to the query domain size, which directly affects the computation and communication costs of the constructions which depend on the database size. In cases when the server database is sparse and only a small fraction of the domain indices correspond to actual entries, using a PIR solution directly will incur a large overhead forcing dependence on the whole domain size. This sparse database setting has been considered as keyword PIR by Chor et al. [15]. The idea of this work is to build an efficiently searchable structure, instantiated with a search tree, over the sparse indices of the database entries and then use PIR to execute the search queries. This approach requires logarithmic number of PIR queries on a database of proportional to the number of sparse items. We propose a new construction which leverages Cuckoo hashing and reduces the overhead to a constant number of PIR queries on a database proportional to the number of data entries.

The idea of our approach is to use Cuckoo hashing to compress the index on the server side. Cuckoo hashing [29, 52] is a dictionary with worst case constant look-up time, which has size linear in the number of inserted items. A Cuckoo hash table is defined by  $\kappa$  hash functions  $H_1, \dots, H_\kappa$  and each item with label  $i$  is placed in one of the  $\kappa$  locations  $H_1(i), \dots, H_\kappa(i)$ . The Cuckoo hash table is initialized by inserting all items in order, resolving collisions using a recursive eviction procedure: whenever an element is hashed to a location that is occupied, the occupying element is evicted and recursively reinserted using a different hash function. For each sequence of items, there is a small set of hash function sets that are incompatible with the sequence and cannot be used to distribute the items, but this can be handled by choosing new hash functions. Overall, inserting  $n$  elements into a cuckoo hash table can be performed in expected  $O(n)$  time [52]. Note that with this procedure the hash functions are dependent on the items placed in the Cuckoo hash table but—unlike in PSI protocols based on Cuckoo hashing [13, 24, 58, 60]—this is not an issue for our use of Cuckoo hashing in the context of PIR where the data is considered public and we do not need to provide any privacy guarantees for it.

Our construction works as follows. The server builds a Cuckoo hash table for its sparse database, which will be of size proportional to the number of non-empty entries (with a constant multiplicative overhead), and provides the Cuckoo hash functions  $H_1, \dots, H_\kappa$  for a  $\kappa \geq 2$ . In order to query an item  $i$ , the client executes  $\kappa$  PIR queries for items  $H_j(i)$ ,  $j \in [\kappa]$  for the database that contains the Cuckoo hash table.

**Table 8: Baseline PIR communication and computation complexities for with different recursion levels and different homomorphic encryption instantiations on a database of size  $n$ .**

PIR protocol	PIR Baseline	PIR Recursion $d = 2$	PIR Recursion $d = \log n$
Additive ElGamal	<b>Comm:</b> $(n + 1) \cdot \lambda_{EG}$ bits <b>Comp:</b> $n$ mults of $\lambda_{EG}$ -bit nums	<b>Comm:</b> $(2n^{1/2} + \lceil F \rceil) \cdot \lambda_{EG}$ bits <b>Comp:</b> $n + n^{1/2} \cdot F$ mults of $\lambda_{EG}$ -bit nums	<b>Comm:</b> $(\log n + \lceil F \rceil^{\log n-1}) \cdot \lambda_{EG}$ bits <b>Comp:</b> $F^{\log n-1}$ mults of $\lambda_{EG}$ -bit nums
Damgård–Jurik (pt = $N^k$ with $N = 2^{\lambda_{DJ}}$ )	<b>Comm:</b> $(n + 1) \cdot (k + 1) \log N$ bits <b>Comp:</b> $n$ mults of $(k + 1)\lambda_{DJ}$ -bit nums	<b>Comm:</b> $n^{1/2}(2k + 3) \log N + (2 + k) \log N$ bits <b>Comp:</b> $n$ mults of $(k + 1)\lambda_{DJ}$ -bit nums + $n^{1/2}$ mults of $(k + 2)\lambda_{DJ}$ -bit nums	<b>Comm:</b> $\approx (k + \log n(1 + k \log n + \log n^2)) \log N$ bits <b>Comp:</b> $n^{i/\log n}$ mults of $(k+1+i)\lambda_{DJ}$ -bit nums for all $i \in \llbracket \log n \rrbracket$ .
Gentry–Ramzan	<b>Comm:</b> $3\lambda_{GR}$ bits <b>Comp:</b> $2 \cdot n \cdot \text{pt}$ multiplications of $\lambda_{GR}$ -bit numbers.	N/A	N/A
FV	<b>Comm:</b> $2(n + 1) \log(q) \cdot \lambda_{FV}$ bits <b>Comp:</b> $n$ scalar mults + additions in $\mathbb{Z}_q[x]/(x^{\lambda_{FV}} + 1)$	<b>Comm:</b> $2(2n^{1/2} + \lceil F \rceil) \log(q) \cdot \lambda_{FV}$ bits <b>Comp:</b> $n + n^{1/2} \lceil F \rceil$ scalar mults + additions in $\mathbb{Z}_q[x]/(x^{\lambda_{FV}} + 1)$	<b>Comm:</b> $2 \log n + \lceil F \rceil^{\log n-1} \log(q) \cdot \lambda_{FV}$ bits <b>Comp:</b> $F^{\log n-1}$ scalar mults + additions in $\mathbb{Z}_q[x]/(x^{\lambda_{FV}} + 1)$

**Table 9: Size of one-time keys required for SealPIR and MulPIR.**

Keys	Recursion	Size (kB)
SealPIR Galois keys	2, 3	6758
MulPIR Galois keys	2	5707
MulPIR Galois keys	3	28270
MulPIR Galois key generator	2	518
MulPIR Galois key generator	3	2356
MulPIR Switching key	2	19
MulPIR Switching key	3	59

We note that our approach to compress the server index using Cuckoo hashing is orthogonal to the use of Cuckoo hashing to batch multiple PIR queries described in Appendix C.3. Combining these two techniques we optimize on two different axis of the PIR construction. Next we present the formal construction for PIR on sparse data.

CONSTRUCTION 1. *Let*

$(\text{Cuckoo.KeyGen}, \text{Cuckoo.Query}, \text{Cuckoo.Insert})$

*be a Cuckoo hash scheme and  $(\text{PIR.Query}, \text{PIR.Eval})$  be a PIR protocol. We construct a new PIR protocol  $(\text{PIR}'.\text{Query}, \text{PIR}'.\text{Eval})$  where the indices of the server's database are sparse over the whole domain:*

- *Pre-processing: The server generates parameters for the Cuckoo hash that will fit its input*

$$(H_1, H_2, \dots, H_\kappa, m) \leftarrow \text{Cuckoo.KeyGen}(|D|).$$

*It initializes the Cuckoo hash table using its input, invoking  $\text{Cuckoo.Insert}(i, d)$  for all  $(i, d) \in D$ . It sends to the client  $\{H_j\}_{j \in [\kappa]}$ .*

- $q_i = (q_i^1, \dots, q_i^\kappa) \leftarrow \text{PIR}'.\text{Query}(i)$ : *The client computes  $q_i^j \leftarrow \text{PIR.Query}(H_j(i))$  for  $j \in [\kappa]$ .*

- $[D[i], \perp] \leftarrow \text{PIR}'.\text{Eval}([q_i, D])$ : *The client and the server run  $[T_j[H_j(i)], \perp] \leftarrow \text{SPIR.Eval}([q_i^j, T_j])$  for  $j \in [\kappa]$ . The client checks if any of the  $T_j[H_j(i)], j \in [\kappa]$  contains item  $i$ . If the items is present, the client outputs it and otherwise, the client outputs  $\perp$ .*

## C.2 Symmetric PIR from OPRFs

The security requirements of a PIR protocol pertain only to the privacy of the query. Symmetric private information retrieval (SPIR) [35] considers also database privacy in addition to query privacy. While some PIR solutions based on homomorphic encryption do effectively provide SPIR guarantees in the case when the server returns a single ciphertext that encrypts only the retrieved database entry, other approaches do provide more information about the database to the client. We provide a simple transformation that enables SPIR given any PIR scheme. Our idea is to encrypt each database entry using a symmetric encryption under a key that is derived in a pseudorandom manner from the index of the data item. In particular, the server derives the encryption keys using pseudorandom function that also offers oblivious evaluation mechanism (OPRF) [30, 41]. To execute a SPIR query the client and the server execute the corresponding PIR query on the database of encrypted entries and in addition to this they run an oblivious PRF evaluation that enables the client to get a single decryption key corresponding to the query entry. We present our protocol next.

CONSTRUCTION 2. *Let  $(\text{Gen}, \text{Enc}, \text{Dec})$  be a semantically secure encryption scheme,  $(\text{PIR.Query}, \text{PIR.Eval})$  be a PIR scheme and  $(\text{PRF.KeyGen}, \text{PRF.Eval}, \text{PIR.OblivEvaluate})$  be an oblivious PRF function. We construct an SPIR protocol as follows:*

- *Pre-processing: The server encrypts its database  $D$  of size  $n$  as follows. It samples a PRF key  $K \leftarrow \text{PRF.KeyGen}(1^\lambda)$  and for each  $i \in [n]$ , it computes  $K_i \leftarrow \text{PRF.Eval}(K, i)$  and sets  $\hat{D}[i] = \text{Enc}(K_i, D[i])$ .*
- $q_i \leftarrow \text{SPIR.Query}(i)$ : *Output  $\text{PIR.Query}(i)$ .*



- $[D[i], \perp] \leftarrow \text{SPIR.Eval}([q_i, D])$ :
  - (1) *The client and the server run  $[\tilde{D}[i], \perp] \leftarrow \text{PIR.Eval}([q_i, \tilde{D}])$ .*
  - (2) *The client and the server evaluate*

$$[K_i, \perp] \leftarrow \text{PRF.OblivEvaluate}([i, K]).$$

- (3) *The client retrieves its output  $D[i] = \text{Dec}(K_i, \tilde{D}[i])$ .*

We note that handling sparse data in the setting of SPIR, requires to use oblivious Cuckoo hashing where the hash function parameters are independent of the data inserted in the hash table. Achieving oblivious Cuckoo hashing requires addition of a stash of size  $O(\log n)$  that stores items which could not be allocated in the hash table due to collisions [44]. The SPIR construction for sparse data proceeds as follows: the server generates a PRF key  $K$  and hash functions for oblivious Cuckoo hashing, it encrypts each item  $i$  in its database with key  $\text{PRF.Eval}(K, i)$ , the server initializes the oblivious Cuckoo hash with the encrypted data. The server sends the Cuckoo hash functions and the encrypted stash to the client. The client executes a query for item  $i$  by running two SPIR queries for  $H_1(i)$  and  $H_2(i)$  using the SPIR construction above. It uses the decryption key  $\text{PRF.Eval}(K, i)$  it has obtained to try to decrypt both the answers in the SPIR queries as well as the encrypted items in the stash. The communication related to the stash can be amortized across different queries.

### C.3 Multi-Query PIR

The traditional definition of PIR considers a setting where queries are executed independently one by one. However, there are scenarios where several queries may be available to be executed at the same time. Multi-query PIR solutions aim to leverage the capability for parallel execution of such queries in order to amortize the complexity. We leverage two main types of techniques for batching: probabilistic batch codes based on Cuckoo hashing [52], which have been used in the context of PIR [5] and private set intersection [25, 61], as well as a CRT batching technique introduced by Groth et al. [39] for Gentry–Ramzan.

### C.4 Private Set + Functionalities

In this section we discuss functionalities which can be solved using specific PIR instantiations. Two such functionalities are private set membership (PSM) and private set intersection (PSI). Private set membership considers the question how to check whether an element held by one party is in the set held by another party. This problem can be viewed as sparse PIR where the database content is the indices themselves. Private set intersection aims to compute the intersection of two private sets. This problem is a generalization of PSM from a single query to multiple queries. Thus, the PSI problem can be phrased as a multi-query PIR on a sparse database. In setting where the two intersection sets have asymmetric sizes, i.e., one of the sets is much smaller, solving PSI using multi-query PIR using the smaller set as queries could provide better asymptotic communication complexity than PSI solutions that require communication linear in the size of the sets.

## D FULL RECURSION WITH LEVELED HOMOMORPHIC ENCRYPTION

In this section, we report on an implementation of *full* recursion  $d = \log n$ , using the technique from Section 4.3. We use the SEAL library [2] with polynomials of degree 8192 and a modulus  $q$  of 147 bits (product of three 49-bit moduli), and plaintext space  $t = 2$ . We implemented full recursion for a Pung-style databases of  $n$  elements of 288B (in particular, we will have one element per ciphertext) [6] and provide benchmarks for databases of size  $2^{14}$  to  $2^{17}$  in Table 10.<sup>4</sup> While this approach does not bring any benefit in practice compared to recursion  $d = 2$  using one homomorphic multiplication, we report for the first time benchmarks for PIR with full recursion.

**Table 10: Full Recursion using Seal-MulPIR.**

$n$	Communication (kB)	Server computation (s)
$2^{14}$	$14 \cdot 150 + 150$	167
$2^{15}$	$15 \cdot 150 + 150$	324
$2^{16}$	$16 \cdot 150 + 150$	658
$2^{17}$	$17 \cdot 150 + 150$	2109

<sup>4</sup>We ran out of RAM for  $n = 2^{18}$  with our tree-based implementation of the tensor product. Careful optimizations of the tensor product computation and regular folding would enable to reduce the memory usage of the program at the cost of increasing computation.