



# One Server for the Price of Two: Simple and Fast Single-Server Private Information Retrieval

Alexandra Henzinger  
MIT

Matthew M. Hong  
MIT

Henry Corrigan-Gibbs  
MIT

Sarah Meiklejohn  
Google

Vinod Vaikuntanathan  
MIT

**Abstract.** We present SimplePIR, the fastest single-server private information retrieval scheme known to date. SimplePIR’s security holds under the learning-with-errors assumption. To answer a client’s query, the SimplePIR server performs fewer than one 32-bit multiplication and one 32-bit addition per database byte. SimplePIR achieves 10 GB/s/core server throughput, which approaches the memory bandwidth of the machine and the performance of the fastest *two-server* private-information-retrieval schemes (which require non-colluding servers). SimplePIR has relatively large communication costs: to make queries to a 1 GB database, the client must download a 121 MB “hint” about the database contents; thereafter, the client may make an unbounded number of queries, each requiring 242 KB of communication. We present a second single-server scheme, DoublePIR, that shrinks the hint to 16 MB at the cost of slightly higher per-query communication (345 KB) and slightly lower throughput (7.4 GB/s/core). Finally, we apply our new private-information-retrieval schemes, together with a novel data structure for approximate set membership, to the task of private auditing in Certificate Transparency. We achieve a strictly stronger notion of privacy than Google Chrome’s current approach with modest communication overheads: 16 MB of download per month, along with 150 bytes per TLS connection.

## 1 Introduction

In a private information retrieval (PIR) protocol [27, 62], a database server holds an array of  $N$  records. A client wants to fetch record  $i \in \{1, \dots, N\}$  from the server, without revealing the index  $i$  that it desires to the server. PIR has applications to systems for private database search [85, 93], metadata-hiding messaging [9, 10], private media consumption [55], credential breach reporting [66, 82, 91, 94], private contact discovery [60], privacy-friendly advertising [11, 54, 58, 87], and private blacklist lookups [61], among others.

Modern PIR schemes require surprisingly little communication: with a single database server and under modest cryptographic assumptions [21, 49, 80], the total communication required to fetch a database record grows only polylogarithmically with the number of records,  $N$ . Unfortunately,

PIR schemes are computationally expensive: the server must touch every bit of the database to answer even a single client query [13], since otherwise the PIR scheme leaks information about which database records the client is *not* interested in. (A number of recent PIR schemes preprocess the database such that the server can answer a query in time sublinear in  $N$ , but all known approaches require either client-specific preprocessing [30, 31, 61, 88, 97] or impractically large server storage [13, 17, 22].) Thus, a hard limit on the throughput of PIR schemes—that is, the ratio between the database size and the server time to answer a query—is the speed with which the PIR server can read the database from memory: roughly 12.4 GB/s/core on our machine [92].

In the standard setting, in which the client interacts with a single database server, the performance of existing PIR protocols is far from this theoretical limit: we measure that the fastest prior single-server PIR schemes [76] achieve a throughput of 259 MB/s/core, or 2% of our machine’s memory bandwidth, on a database of hundred-byte records. It is possible to push the performance up to 1.3 GB/s/core when the database records are hundreds of kilobytes long, though that parameter setting is not relevant for many PIR applications, including our application to Certificate Transparency.

When the client can communicate with multiple *non-colluding* database servers [27], there exist PIR schemes with server-side throughput of up to 11.5 GB/s/core, or 93% of the memory bandwidth (described in Table 1). However, these multi-server PIR schemes are cumbersome to deploy, since they rely on multiple coordinating yet independent infrastructure providers. In addition, their security is brittle, as it stems from a non-collusion assumption rather than from cryptographic hardness. Thus, existing PIR schemes suffer from either poor performance—in the single-server setting—or undesirable trust assumptions—in the multi-server case.

In this paper, we present two new single-server PIR schemes that exceed the throughput of all existing single-server PIR protocols and approach the throughput of multi-server ones. In addition, our schemes are relatively simple to explain and easy to implement: our complete implementation of both schemes, available at [github.com/ahenzinger/simplepir](https://github.com/ahenzinger/simplepir), requires roughly 1,400 lines of Go code, plus 200 lines of C, and uses no external libraries.

More specifically, our first scheme, SimplePIR, achieves a server throughput of 10 GB/s/core, or 81% of the memory

This is the full version of a paper of the same title at USENIX Security 2023.

bandwidth, though it requires the client to download a relatively large “hint” about the database contents before making its queries. On a database of  $N$  bytes, the hint has size roughly  $4\sqrt{N}$  KB. The hint is not client-specific, and a client can reuse the hint over many queries, so the amortized communication cost per query can be small. Our second scheme, DoublePIR, achieves slightly lower server throughput of 7.4 GB/s/core, but shrinks the hint to roughly 16 MB for a database of one-byte records—independent of the number of records in the database.

**Our techniques.** We now summarize the technical ideas behind our results.

*Recap: Single-server PIR.* Our starting point is the single-server PIR construction of Kushilevitz and Ostrovsky [62]. In their scheme, the PIR server represents an  $N$ -record database as a matrix  $\mathbf{D}$  of dimension  $\sqrt{N}$  by  $\sqrt{N}$ . To fetch the database record in row  $i$  and column  $j$ , the client sends the server the encryption  $E(\mathbf{q})$  of a dimension- $\sqrt{N}$  vector that is zero everywhere except that it has a “1” in index  $j$ . If the encryption scheme is linearly homomorphic, the server can compute the matrix-vector product  $\mathbf{D} \cdot E(\mathbf{q}) = E(\mathbf{D} \cdot \mathbf{q})$  under encryption and return the result to the client. The client decrypts to recover  $\mathbf{D} \cdot \mathbf{q}$  which, by construction, is the  $j$ -th column of the database, as desired. The total communication grows as  $\sqrt{N}$ .

*SimplePIR from linearly homomorphic encryption with preprocessing.* The PIR server’s throughput here is limited by the speed with which it can compute the product of the plaintext matrix  $\mathbf{D}$  with the encrypted vector  $E(\mathbf{q})$ . Our observation in SimplePIR (Section 4) is that, using Regev’s learning-with-errors-based encryption scheme [86], the server can perform the vast majority of the work of computing the matrix-vector product  $\mathbf{D} \cdot E(\mathbf{q})$  in advance—before the client even makes its query. The server’s preprocessing depends only on the database  $\mathbf{D}$  and the public parameters of the Regev encryption scheme, so the server can reuse this preprocessing work across many queries from many independent clients. After this preprocessing step, to answer a client’s query, the server needs to compute only roughly  $N$  32-bit integer multiplications and additions on a database of  $N$  bytes. The catch is that the client must download a “hint” about the database contents after this preprocessing step—the hint accounts for the bulk of the communication cost in SimplePIR.

*DoublePIR from one recursive step.* The idea behind DoublePIR (Section 5) comes from the original Kushilevitz and Ostrovsky paper [62]: in SimplePIR, the client downloads the hint from the server, along with a dimension- $\sqrt{N}$  encrypted vector. However, to recover its record of interest, the client only needs one small part of the hint and one component of this vector. We show how the client can use SimplePIR recursively on the hint and this vector to fetch its desired database record at a reduced communication cost. To minimize the concrete costs, we make non-black-box use of SimplePIR in this recursive construction, which saves a factor of the lattice dimension, which is 1024 for our parameters, over a naïve design.

**Application to Certificate Transparency.** Finally, we evaluate our PIR schemes in the context of the application of *signed certificate timestamp (SCT) auditing* in Certificate Transparency. In this auditing application, a server holds a set  $S$  of strings and a client (web browser) wants to test whether a particular string  $\sigma$ , representing an SCT, appears in the set  $S$ , while hiding  $\sigma$  from the server. (The string  $\sigma$  reveals information about which websites a client has visited.) Google Chrome currently implements this auditing step using a solution that provides  $k$ -anonymity for  $k = 1000$  [35].

Along the way, we construct a new data structure (Section 6) for more efficiently solving this type of private set-membership problem using PIR, when a constant rate of false positives is acceptable (as in our application). In this setting, standard Bloom filters [15] and approaches based on PIR by keywords [26] require the client to perform PIR over a database of  $\lambda N$  bits (if the set  $S$  has size  $N$  and  $\lambda \approx 128$  is a security parameter). In contrast, our data structure requires performing PIR over only  $8N$  bits—giving a roughly 16 $\times$  speedup in our application.

Google’s current solution to SCT auditing, which provides  $k$ -anonymity rather than full cryptographic privacy, requires the client to communicate 24 B on average per TLS connection. Our solution, which provides cryptographic privacy, requires 150 B and 0.0003 core-seconds of server compute on average per TLS connection, along with 16 MB of client download and 150 KB of client storage every month to maintain the hint.

**Limitations.** Our new PIR schemes come with two main downsides. First, our client must download a “hint”: on databases gigabytes in size, the hint is tens of megabytes. If a client makes only one query, this hint download dominates the overall communication. Second, our schemes’ online communication is on the order of hundreds of kilobytes, which is 10 $\times$  larger than in some prior work. Nevertheless, we believe that SimplePIR and DoublePIR represent an exciting new point in the PIR design space: large computation savings, along with a conceptually simple design and small, stand-alone codebase, at the cost of modest communication and storage overheads.

**Our contributions.** In summary, our contributions are:

- two new high-throughput single-server private information retrieval protocols (Sections 4 and 5),
- a new data structure for private set membership using PIR (Section 6) and its application to private auditing in Certificate Transparency (Section 7), and
- the evaluation of these schemes, using a new open-source implementation (Section 8).

## 2 Related work and comparison

Chor, Goldreich, Kushilevitz and Sudan [27] introduced PIR in the multi-server setting and Kushilevitz and Ostrovsky [62] gave the first construction of single-server PIR. Their scheme uses a linearly homomorphic encryption scheme that expands

Scheme	Servers	Communication	No per-client storage (on the server)	Polylog( $n$ ) overhead	Max. achievable throughput/core
DPF PIR [16, 59]	2	$\log N$	✓	✓	5,381 MB/s
XOR PIR [27]	2	$\sqrt{N}$	✓	✓	6,067 MB/s
XOR PIR fast* [27]	2	$\sqrt{N}$	✓	✓	11,797 MB/s
SealPIR [9] ( $d = 2$ )	1	$\sqrt{N}$	✗	✓	97 MB/s
MulPIR [8] ( $d = 2$ )	1	$\sqrt{N}$	✗	✓	69 MB/s <sup>◇</sup>
FastPIR [6]	1	$N$	✗	✓	215 MB/s
OnionPIR [77]	1	$\log N$	✗	✓	104 MB/s
Spiral family [76]	1	$\log N$	✗	✓	1,314 MB/s
KO [62]+Paillier [81]	1	$N^\epsilon$	✓	✗	0.131 MB/s
XPIR [5] ( $d = 2$ )	1	$\sqrt{N}$	✓	✓	142 MB/s <sup>◇</sup>
FrodoPIR* [34]	1	$\sqrt{N}$	✓	✓	1,256 MB/s
<b>SimplePIR</b> (§4)	1	$\sqrt{N}$	✓	✓	10,305 MB/s
<b>DoublePIR</b> (§5)	1	$\sqrt{N}$	✓	✓	7,622 MB/s

Table 1: A comparison of PIR schemes on database size  $N$  and security parameter  $n$ . The overhead column indicates whether the server computation per database bit is at most polylogarithmic in  $n$ . The throughput column gives the maximum throughput we measured for any record size. The database and record sizes used are in Appendix A. The throughput is normalized by the number of cores, i.e., divided by two for two-server schemes. \*This is a non-constant-time implementation—each server’s running time depends on its secret input. We include the performance for comparison, though a side-channel-resistant production implementation might not use this optimization. <sup>◇</sup>No open-source code available; this throughput is reported in the MulPIR paper [8]. <sup>◇</sup>This XPIR throughput is reported by SealPIR [9]. \*FrodoPIR is concurrent work and is essentially identical to SimplePIR, up to the choice of lattice parameters (see Section 2).

$\ell$ -bit plaintexts to  $\ell \cdot F$ -bit ciphertexts. We call  $F$  the *expansion factor* of the encryption scheme. Then, on a database of  $N$  bits and any dimension parameter  $d \in \{1, 2, 3, \dots\}$ , their PIR construction has communication roughly  $N^{1/d} F^{d-1}$ . The server must perform roughly  $NF^{d-1}$  homomorphic operations in the process of answering the client’s query.

The Damgård-Jurik [33] cryptosystem has expansion factor  $F \approx 1 + \epsilon$  for any constant  $\epsilon > 0$ , which yields very communication-efficient PIR schemes [69]. It is possible to construct PIR with similar communication efficiency from an array of cryptographic assumptions [21, 24, 39]. However, these schemes are all costly in computation: for each *bit* of the database, the server must perform work polynomial in the security parameter.

**Lattice-based PIR.** To drive down this computational cost, recent PIR schemes instantiate the Kushilevitz-Ostrovsky construction using encryption schemes based on the ring learning-with-errors problem (“Ring LWE”) [72]. In these schemes, for each bit of the database, the server performs work *polylogarithmic* in the security parameter—rather than polynomial.

However, these savings in computation come at the cost of a larger expansion factor ( $F \approx 10$ ), which increases the communication as the dimension parameter  $d$  cannot be too large. For example, XPIR [5] takes  $d = 2$ . In addition, the client in the Kushilevitz-Ostrovsky scheme must upload  $N^{1/d}$  ciphertexts, and each ring-LWE ciphertext is at least thousands of kilobytes in size. This imposes large absolute communication costs (e.g., tens of MB per query, on a database of hundreds of MB).

SealPIR [9] shows that the client can *compress* the ciphertexts in an XPIR-style scheme before uploading them. The server can then expand these ciphertexts using homomorphic operations. (FastPIR [6] uses a similar idea to compress responses.) This optimization reduces the communication costs by orders of magnitude, though it requires the server to store some per-client information (“key-switching hints”)—essentially, encryptions of the client’s secret decryption key—that is megabytes in size and that the client must upload to the server before it makes any queries.

MulPIR [8], OnionPIR [77], and Spiral [76] additionally use *fully* homomorphic encryption [47] to reduce the communication cost. In Spiral [76], for example, the cost grows roughly as  $N^{1/d} F$ , where the exponent on the  $F$  term is now 1 instead of  $d - 1$ . Building on ideas of Gentry and Halevi [48], Spiral shows how to decrease the communication cost while keeping the throughput high: up to 259 MB/s on a database of short records. (With long database records, Spiral does not use the SealPIR query compression technique and gets throughput as large as 1,314 MB/s, at the cost of increased communication.)

**Plain learning with errors.** We base our PIR schemes on the standard learning-with-errors (LWE) problem—not the ring variant. The expansion factor of the standard LWE-based encryption scheme, Regev encryption [86], is roughly  $F = n \approx 1024$ , where  $n$  is the lattice security parameter. This large expansion factor means that a direct application of Regev encryption to the Kushilevitz-Ostrovsky PIR scheme would be disastrous in terms of communication and computation. Our innovation is to show that the server can do the bulk of its work *in advance*, and reuse it over multiple clients.

Aside from the fact that our scheme is based on a *weaker cryptographic assumption*, namely plain LWE as opposed to ring LWE, this strategy yields multiple benefits:

1. Our LWE-based schemes are *simple* to implement: they require no polynomial arithmetic or fast Fourier transforms.
2. Our schemes do not require the server to store any extra per-client state. In contrast, many schemes based on Ring LWE [8, 9, 76, 77] rely on optimizations that require the server to store one “key-switching hint” for each client.
3. Our schemes are *faster*. We avoid the costs associated with ciphertext compression and expansion. In addition, since we only need our encryption scheme to be linearly (*not* fully) homomorphic, we can use smaller and more efficient lattice parameters.

The drawback of our schemes is that they have larger com-

munication cost, especially when the client makes only a single query (so the client cannot amortize the offline download cost over multiple queries) or when the database records are long.

**Concurrent work: FrodoPIR.** FrodoPIR [34] is independent concurrent work that constructs a PIR scheme that is essentially identical to SimplePIR. The default configuration of FrodoPIR has communication cost  $O(N)$ , on database size  $N$ , though rebalancing the scheme gives a  $O(n\sqrt{N})$ -cost, on lattice dimension  $n$ , as in SimplePIR. The additional contributions of our work are: the more communication-efficient DoublePIR scheme, our new data structure for private set-membership queries (Section 6), the application to certificate transparency (Section 7), and an optimized implementation of our schemes.

**Preprocessing and PIR.** The server in our PIR schemes performs some client-independent preprocessing. Prior work uses server-side preprocessing—either one-time [13, 17, 22] or per-client [30, 31, 61, 88, 97]—to build PIR where the server online work is *sublinear* in the database size. Prior work also proves strong lower bounds on the performance of any such PIR with preprocessing schemes [13, 30, 31, 84]. In contrast, in this work, we use preprocessing to build PIR where the amortized per-query server work is still *linear*, but it is concretely efficient.

**Multi-server PIR.** In our PIR schemes, the client communicates with a single database server. In multi-server PIR schemes [27], the client communicates with multiple database servers and client privacy holds only as long an attacker cannot compromise some number of them. In Table 1, we give the throughput of an optimized implementation [59] of a two-server PIR scheme based on distributed point functions (“DPF PIR”) [16, 50]. We also report the throughput of a  $\sqrt{N}$ -communication two-server PIR scheme (“XOR PIR”) [27]. It is possible to speed these schemes up by roughly  $2\times$  if the server’s running time can depend on the *Hamming weight* of the client’s query vector (“XOR PIR fast”). The downside of this optimization is that it could potentially leak information about one server’s secret query vector to another server via timing information, thereby breaking client privacy. Whether such a performance-leakage trade-off is acceptable in practice likely depends on the application scenario.

**Hardware acceleration for PIR.** Recent work improves the throughput of both single-server [67] and multi-server [56] PIR using hardware acceleration. This approach is complementary to ours, as it may further speed up our new PIR protocols.

**Privacy and certificate transparency.** Lueks and Goldberg [71] and Kales, Omolola, and Ramacher [59] propose using *multi-server* PIR for auditing in certificate transparency. We work in the *single-server* setting, where the client communicates with a separate audit server (e.g., Google, in the application to Chrome). Further, we introduce a new set-membership data structure to reduce the cost of auditing (Section 6). We discuss existing approaches to auditing in Section 7.

## 3 Background and definitions

**Notation.** For a probability distribution  $\chi$ , we use  $x \leftarrow^{\mathbb{R}} \chi$  to indicate that  $x$  is a random sample from  $\chi$ . For a finite set  $S$ , we use  $x \leftarrow^{\mathbb{R}} S$  to denote sampling  $x$  uniformly at random from  $S$ . We use  $\mathbb{N}$  to represent the natural numbers and  $\mathbb{Z}_p$  to represent integers modulo  $p$ . All logarithms are to the base two. For  $x \in \mathbb{N}$ , we let  $[x]$  denote the set  $\{1, \dots, x\}$ . Throughout, we assume that values like  $\sqrt{N}$  are integral, wherever doing so is essentially without loss of generality. Algorithms are modeled as RAM programs and their runtime is measured in terms of the number of RAM instructions executed. We use the symbols MB and GB to denote  $2^{20}$  and  $2^{30}$  bytes, respectively.

### 3.1 Learning with errors (LWE)

The security of our PIR schemes relies on the decision version of the learning-with-errors assumption [86]. The assumption is parameterized by the dimension of the LWE secret  $n \in \mathbb{N}$ , the number of samples  $m \in \mathbb{N}$ , the integer modulus  $q \geq 2$ , and an error distribution  $\chi$  over  $\mathbb{Z}$ . The LWE assumption then asserts that for a matrix  $\mathbf{A} \leftarrow^{\mathbb{R}} \mathbb{Z}_q^{m \times n}$ , a secret  $\mathbf{s} \leftarrow^{\mathbb{R}} \mathbb{Z}_q^n$ , an error vector  $\mathbf{e} \leftarrow^{\mathbb{R}} \chi^m$ , and a random vector  $\mathbf{r} \leftarrow^{\mathbb{R}} \mathbb{Z}_q^m$ , the following distributions are computationally indistinguishable:

$$\{(\mathbf{A}, \mathbf{A}\mathbf{s} + \mathbf{e})\} \stackrel{c}{\approx} \{(\mathbf{A}, \mathbf{r})\}.$$

More specifically, the  $(n, q, \chi)$ -LWE problem with  $m$  samples is  $(T, \epsilon)$ -hard if all adversaries running in time  $T$  have advantage at most  $\epsilon$  in distinguishing the two distributions. In Section 4.2, we give concrete values for the LWE parameters.

**Secret-key Regev encryption.** Regev [86] gives a secret-key encryption scheme that is secure under the LWE assumption. With LWE parameters  $(n, q, \chi)$  and a plaintext modulus  $p$ , the Regev secret key is a vector  $\mathbf{s} \leftarrow^{\mathbb{R}} \mathbb{Z}_q^n$ . The Regev encryption of a message  $\mu \in \mathbb{Z}_p$  is

$$(\mathbf{a}, c) = (\mathbf{a}, \mathbf{a}^\top \mathbf{s} + e + \lfloor q/p \rfloor \cdot \mu) \in \mathbb{Z}_q^n \times \mathbb{Z}_q,$$

for  $e \leftarrow^{\mathbb{R}} \chi$ . To decrypt the ciphertext, anyone who knows the secret  $\mathbf{s}$  can compute  $c - \mathbf{a}^\top \mathbf{s} \bmod q$  and round the result to the nearest multiple of  $\lfloor q/p \rfloor$ . Decryption succeeds as long as the absolute value of the error sampled from the error distribution  $\chi$  is smaller than  $\frac{1}{2} \cdot \lfloor q/p \rfloor$ . We say that a setting of the Regev parameters supports *correctness error*  $\delta$  if the probability of a decryption error is at most  $\delta$  (over the encryption algorithm’s randomness). Regev encryption is additively homomorphic, since given two ciphertexts  $(\mathbf{a}_1, c_1)$  and  $(\mathbf{a}_2, c_2)$ , their sum  $(\mathbf{a}_1 + \mathbf{a}_2, c_1 + c_2)$  decrypts to the sum of the plaintexts, provided again that the error remains sufficiently small.

### 3.2 Private information retrieval with hints

We now give the syntax and security definitions for the type of PIR schemes we construct. Our form of PIR is very similar to the standard single-server PIR schemes [27, 62]. The primary

distinction is that we allow the PIR server to preprocess the database ahead of time and to output two “hints”: one that the server stores locally, and another that the server sends to each client. This preprocessing allows the PIR server to push much of its computational work into an offline phase that takes place before the client makes its query. In our constructions, both hints are small—they have size sublinear in the database size. In addition, all clients use the same hint and a client can reuse the same hint for all of its of PIR queries.

*Remark 3.1 (Handling database updates).* As PIR schemes with preprocessing perform some precomputation over the database, the server inherently needs to repeat some of this work if the database contents change. Related work investigates how to minimize the amount of computation and communication that such database updates incur, in both a black-box [61] and a protocol-specific [73] manner. We address how to handle updates in our schemes in Appendices C.3 and E.3.

A PIR-with-preprocessing scheme [13], over plaintext space  $\mathcal{D}$  and database size  $N \in \mathbb{N}$ , consists of four routines, which all take the security parameter as an implicit input:

Setup(db)  $\rightarrow$  (hint<sub>s</sub>, hint<sub>c</sub>). Given a database db  $\in \mathcal{D}^N$ , output preprocessed hints for the server and the client.

Query(*i*)  $\rightarrow$  (st, qu). Given an index *i*  $\in [N]$ , output a secret client state st and a database query qu.

Answer(db, hint<sub>s</sub>, qu)  $\rightarrow$  ans. Given the database db, a server hint hint<sub>s</sub>, and a client query qu, output an answer ans.

Recover(st, hint<sub>c</sub>, ans)  $\rightarrow$  *d*. Given a secret client state st, a client hint hint<sub>c</sub>, and an answer ans, output a record *d*  $\in \mathcal{D}$ .

**Correctness.** When the client and the server execute the PIR protocol faithfully, the client should recover its desired database record with all but negligible probability in the implicit correctness parameter. Formally, we say that a PIR scheme has *correctness error*  $\delta$  if, on database size  $N \in \mathbb{N}$ , for all databases db = (*d*<sub>1</sub>, . . . , *d*<sub>*N*</sub>)  $\in \mathcal{D}^N$  and for all indices *i*  $\in [N]$ , the following probability is at least  $1 - \delta$ :

$$\Pr \left[ d_i = \hat{d}_i : \begin{array}{l} (\text{hint}_s, \text{hint}_c) \leftarrow \text{Setup}(\text{db}) \\ (\text{st}, \text{qu}) \leftarrow \text{Query}(i) \\ \text{ans} \leftarrow \text{Answer}(\text{db}, \text{hint}_s, \text{qu}) \\ \hat{d}_i \leftarrow \text{Recover}(\text{st}, \text{hint}_c, \text{ans}) \end{array} \right].$$

For the PIR scheme to be non-trivial, the total client-to-server communication should be smaller than the bitlength of the database. That is, it must hold that  $|\text{hint}_c| + |\text{qu}| + |\text{ans}| \ll |\text{db}|$ .

**Security.** The client’s query should reveal no information about its desired database record. That is, we say that a PIR scheme is (*T*,  $\epsilon$ )-*secure* if, for all adversaries  $\mathcal{A}$  running in time at most *T*, on database size  $N \in \mathbb{N}$ , and for all *i*, *j*  $\in [N]$ ,

$$\left| \Pr[\mathcal{A}(1^N, \text{qu}) = 1 : (\text{st}, \text{qu}) \leftarrow \text{Query}(i)] - \Pr[\mathcal{A}(1^N, \text{qu}) = 1 : (\text{st}, \text{qu}) \leftarrow \text{Query}(j)] \right| \leq \epsilon.$$

*Remark 3.2 (Stateless client).* The client in our PIR schemes does not hold any secret state across queries. In contrast, in SealPIR [9] and related schemes, the client builds its queries using persistent, long-term cryptographic secrets. We show in Appendix B that, in certain settings, a malicious PIR server can perform a state-recovery attack against these schemes and thus break client privacy for both past and future queries. Our stateless schemes are not vulnerable to such attacks.

## 4 SimplePIR

In this section, we present our first PIR scheme, SimplePIR. SimplePIR is the fastest single-server PIR scheme known to date in terms of throughput per second per core (Table 1). In particular, we prove the following theorem:

**Informal Theorem 4.1.** *On database size  $N$ , let  $p \in \mathbb{N}$  be a suitable plaintext modulus for secret-key Regev encryption with LWE parameters  $(n, q, \chi)$ , achieving  $(T, \epsilon)$ -security for  $\sqrt{N}$  LWE samples and supporting  $\sqrt{N}$  homomorphic additions with correctness error  $\delta$  (cf. Section 4.2). Then, for a random LWE matrix  $\mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n}$ , SimplePIR is a  $(T - O(\sqrt{N}), 2\epsilon)$ -secure PIR scheme on database size  $N$ , over plaintext space  $\mathbb{Z}_p$ , with correctness error  $\delta$ .*

We give a formal description of SimplePIR in Figure 2 and we prove its security and correctness in Appendix C.

*Remark 4.1 (Concrete costs of SimplePIR).* Using the parameters of Informal Theorem 4.1, we give SimplePIR’s concrete costs, with no hidden constants, in terms of operations (i.e., integer additions and multiplications) over  $\mathbb{Z}_q$ . In a one-time public preprocessing phase, SimplePIR requires:

- the server to perform  $2nN$  operations in  $\mathbb{Z}_q$ , and
- the client to download  $n\sqrt{N}$  elements in  $\mathbb{Z}_q$ ,

where our implementation takes  $n = 2^{10}$  and  $q = 2^{32}$  to achieve 128-bit security against the best known attacks [7].

On each query, SimplePIR requires

- the client to upload  $\sqrt{N}$  elements in  $\mathbb{Z}_q$ ,
- the server to perform  $2N$  operations in  $\mathbb{Z}_q$ , and
- the client to download  $\sqrt{N}$  elements in  $\mathbb{Z}_q$ .

### 4.1 Technical ideas

We now discuss the SimplePIR construction in more detail.

The simplest non-trivial single-server PIR schemes [23, 62, 69] take the following “square-root” approach: given an *N*-element database, the server stores this database as a  $\sqrt{N}$ -by- $\sqrt{N}$  square matrix. Meanwhile, a client who wishes to query for database entry *i*  $\in [N]$  decomposes index *i* into the pair of coordinates  $(i_{\text{row}}, i_{\text{col}}) \in [\sqrt{N}]^2$ . Then, the client builds a unit vector  $\mathbf{u}_{i_{\text{col}}}$  in  $\mathbb{Z}_2^{\sqrt{N}}$  (i.e., the vector of all zeros with a single ‘1’ at index *i*<sub>col</sub>), element-wise encrypts it with a linearly homomorphic encryption scheme, and sends this encrypted

**Construction: SimplePIR.** The parameters of the construction are a database size  $N$ , LWE parameters  $(n, q, \chi)$ , a plaintext modulus  $p \ll q$ , and a LWE matrix  $\mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n}$  (sampled in practice using a hash function). The database consists of  $N$  values in  $\mathbb{Z}_p$ , which we represent as a matrix in  $\mathbb{Z}_p^{\sqrt{N} \times \sqrt{N}}$ . Define the scalar  $\Delta := \lfloor q/p \rfloor \in \mathbb{Z}$ .

Setup( $\text{db} \in \mathbb{Z}_p^{\sqrt{N} \times \sqrt{N}}$ )  $\rightarrow$  ( $\text{hint}_s, \text{hint}_c$ ).

- Return  $(\text{hint}_s, \text{hint}_c) \leftarrow (\perp, \text{db} \cdot \mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n})$ .

Query( $i \in [N]$ )  $\rightarrow$  ( $\text{st}, \text{qu}$ ).

- Write  $i$  as a pair  $(i_{\text{row}}, i_{\text{col}}) \in [\sqrt{N}]^2$ .
- Sample  $\mathbf{s} \stackrel{\mathcal{R}}{\leftarrow} \mathbb{Z}_q^n$  and  $\mathbf{e} \stackrel{\mathcal{R}}{\leftarrow} \chi^{\sqrt{N}}$ .
- Compute  $\text{qu} \leftarrow (\mathbf{A}\mathbf{s} + \mathbf{e} + \Delta \cdot \mathbf{u}_{i_{\text{col}}}) \in \mathbb{Z}_q^{\sqrt{N}}$ , where  $\mathbf{u}_{i_{\text{col}}}$  is the vector of all zeros with a single ‘1’ at index  $i_{\text{col}}$ .
- Return  $(\text{st}, \text{qu}) \leftarrow ((i_{\text{row}}, \mathbf{s}), \text{qu})$ .

Answer( $\text{db} \in \mathbb{Z}_p^{\sqrt{N} \times \sqrt{N}}, \text{hint}_s, \text{qu} \in \mathbb{Z}_q^{\sqrt{N}}$ )  $\rightarrow$  ans.

- Return  $\text{ans} \leftarrow \text{db} \cdot \text{qu} \in \mathbb{Z}_q^{\sqrt{N}}$ .

Recover( $\text{st}, \text{hint}_c \in \mathbb{Z}_q^{\sqrt{N} \times n}, \text{ans} \in \mathbb{Z}_q^{\sqrt{N}}$ )  $\rightarrow$   $d$ .

- Parse  $(i_{\text{row}} \in [\sqrt{N}], \mathbf{s} \in \mathbb{Z}_q^n) \leftarrow \text{st}$ .
- Compute  $\hat{d} \leftarrow (\text{ans}[i_{\text{row}}] - \text{hint}_c[i_{\text{row}}, :] \cdot \mathbf{s}) \in \mathbb{Z}_q$ , where  $\text{ans}[i_{\text{row}}]$  denotes component  $i_{\text{row}}$  of  $\text{ans}$  and  $\text{hint}_c[i_{\text{row}}, :]$  denotes row  $i_{\text{row}}$  of  $\text{hint}_c$ .
- Return  $d \leftarrow \text{Round}_\Delta(\hat{d})/\Delta \in \mathbb{Z}_p$ , which is  $\hat{d}$  rounded to the nearest multiple of  $\Delta$  and then divided by  $\Delta$ .

Figure 2: The SimplePIR protocol.

vector to the server. The server computes the matrix-vector product between the database and the query vector and returns it to the client. Finally, the client decrypts element  $i_{\text{row}}$  of the server’s answer vector—which corresponds exactly to the inner product of database row  $i_{\text{row}}$  and encrypted unit vector  $\mathbf{u}_{i_{\text{col}}}$ , or, equivalently, the encrypted database entry at  $(i_{\text{row}}, i_{\text{col}})$ . In this scheme, the server and the client exchange  $2\sqrt{N}$  ciphertext elements, while the server performs  $N$  ciphertext multiplications and additions to answer each PIR query.

Our starting point is to instantiate this “square-root” approach with the secret-key version of Regev’s LWE-based encryption scheme [86]. Let  $(n, q, \chi)$  be LWE parameters. Then, the Regev encryption of a vector  $\boldsymbol{\mu} \in \mathbb{Z}_p^m$  consists of a pair of a matrix and a vector:

$$\text{Enc}(\boldsymbol{\mu}) = (\mathbf{A}, \mathbf{c}) = (\mathbf{A}, \mathbf{A}\mathbf{s} + \mathbf{e} + \lfloor q/p \rfloor \cdot \boldsymbol{\mu}),$$

for some LWE matrix  $\mathbf{A} \stackrel{\mathcal{R}}{\leftarrow} \mathbb{Z}_q^{m \times n}$ , secret  $\mathbf{s} \stackrel{\mathcal{R}}{\leftarrow} \mathbb{Z}_q^n$ , and error vector  $\mathbf{e} \stackrel{\mathcal{R}}{\leftarrow} \chi^m$ .

We make three crucial observations about Regev encryption:

1. First, a large part of the ciphertext—namely, the matrix

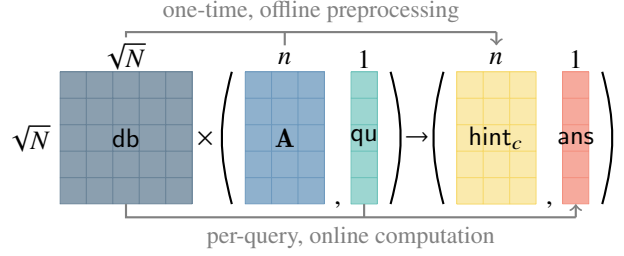


Figure 3: The server computation in SimplePIR. Each cell represents a  $\mathbb{Z}_q$  element, and  $\times$  denotes matrix multiplication. The server performs the bulk of its work in a one-time preprocessing step. Thereafter, the server can answer each client’s query with a lightweight online phase.

$\mathbf{A}$ —is independent of the encrypted message. It is thus possible to generate the matrix  $\mathbf{A}$  ahead of time.

2. Second, Regev encryption remains secure even when the same matrix  $\mathbf{A}$  is used to encrypt polynomially many messages (cf. Corollary C.3), provided that each ciphertext uses an independent secret vector  $\mathbf{s}$  and error vector  $\mathbf{e}$  [83].
3. Finally, we can take  $\mathbf{A}$  to be pseudorandom (rather than random) at a negligible loss in security, allowing us to succinctly represent  $\mathbf{A}$  by a short random seed.

In SimplePIR, we leverage these three observations as follows. Consider a client who wishes to retrieve the database entry at  $(i_{\text{row}}, i_{\text{col}})$ . At a conceptual level, the client’s query to the server consists of  $\text{Enc}(\mathbf{u}_{i_{\text{col}}}) = (\mathbf{A}, \mathbf{c})$ —the Regev encryption of the vector in  $\mathbb{Z}_p^{\sqrt{N}}$  that is zero everywhere but with a “1” at index  $i_{\text{col}}$ . The server then represents the database as a matrix  $\mathbf{D} \in \mathbb{Z}_p^{\sqrt{N} \times \sqrt{N}}$  and computes and returns the matrix-vector product of the database with the client’s encrypted query, i.e.,  $(\mathbf{D} \cdot \mathbf{A}, \mathbf{D} \cdot \mathbf{c})$ . From the server’s reply, the client can use standard Regev decryption to recover  $\mathbf{D} \cdot \mathbf{u}_{i_{\text{col}}} \in \mathbb{Z}_p^{\sqrt{N}}$ , which is exactly the  $i_{\text{col}}$ -th column of the database, as desired.

Now, we make the following modifications:

1. We have the server compute the value  $\mathbf{D} \cdot \mathbf{A}$  ahead of time in a preprocessing phase. This preprocessing step requires  $2nN$  operations in  $\mathbb{Z}_q$ , on lattice dimension  $n \approx 2^{10}$  and database size  $N$ . Then, to answer the client’s query, the server needs to compute the value  $\mathbf{D} \cdot \mathbf{c}$ , which requires only  $2N$  operations in  $\mathbb{Z}_q$ . So, an  $n/(n+1)$  fraction (i.e., 99.9%) of the server’s work can happen ahead of time—before the client even decides which database record it wants to fetch.
2. We have all clients use the same matrix  $\mathbf{A}$  to build each of their queries. Then, the server only precomputes  $\mathbf{D} \cdot \mathbf{A}$  once. The server sends this one-time “hint” to all clients. Thus, the server amortizes the cost of computing and communicating  $\mathbf{D} \cdot \mathbf{A}$  over many clients and over many queries.
3. As an optimization, we compress  $\mathbf{A}$  using pseudorandomness. Specifically, the server and the clients can derive  $\mathbf{A}$  as the output of a public hash function, modelled as a random oracle, applied to a fixed string in counter mode. This saves on bandwidth and storage, as the server and the clients

communicate and store only a small seed to generate  $\mathbf{A}$ .

The security of the SimplePIR construction follows almost immediately from the security of Regev encryption [86] with a reused matrix  $\mathbf{A}$  [83], which in turn follows from the hardness of LWE. SimplePIR’s correctness follows from the correctness of Regev’s linearly homomorphic encryption scheme and of Kushilevitz and Ostrovsky’s “square-root” PIR template.

## 4.2 Parameter selection

Picking the LWE parameters  $(n, q, \chi)$  and the plaintext modulus  $p$  requires a standard (though tedious) analysis. We choose our parameters to have 128-bit security, according to modern lattice-attack-cost estimates [7]. In particular, we set the secret dimension  $n = 2^{10}$ , use modulus  $q = 2^{32}$  (as modern hardware natively supports operations with this modulus), set the error distribution  $\chi$  to be the discrete Gaussian distribution with standard deviation  $\sigma = 6.4$ , and allow correctness error  $\delta = 2^{-40}$ . We obtain the following trade-off between database size  $N$  and plaintext modulus  $p$ :

Database size $N$ :	$2^{26}$	$2^{28}$	$2^{30}$	$2^{34}$	$2^{38}$	$2^{42}$
Plaintext modulus $p$ :	991	833	701	495	350	247

We discuss parameter selection further in Appendix C.1.

## 4.3 Extensions

Finally, we extend our SimplePIR construction to meet the requirements of realistic deployment scenarios:

**Supporting databases with larger record sizes.** The basic SimplePIR scheme (Figure 2) supports a database in which each record is a single  $\mathbb{Z}_p$  element—or, roughly 8-10 bits with our parameter settings. Our main application (Section 7) uses a database with one-bit records, though other applications of PIR [6, 9, 10, 55, 76] use much longer records.

To handle large records, we observe that the client in SimplePIR can retrieve an entire column of the database at once. Concretely, after executing a single online phase with the server to query for database element  $(i_{\text{row}}, i_{\text{col}})$ , the client can run the Recover procedure  $\sqrt{N}$  times—once for every row in  $[\sqrt{N}]$ —to reconstruct the entire column  $i_{\text{col}}$  of the database matrix. So, to support large records, we encode each record as multiple elements in the plaintext space,  $\mathbb{Z}_p$ , and store these elements stacked vertically in the same column. By making a single online query and reconstructing the corresponding column of elements, the client recovers any record of its choosing.

On a database of  $N$  records, each in  $\mathbb{Z}_p^d$  (where  $d \leq N$ ), with LWE secret dimension  $n$  and modulus  $q$ , SimplePIR has:

- one-time (hint) download  $n \cdot \sqrt{dN}$  elements in  $\mathbb{Z}_q$ ,
- per-query upload and download  $\sqrt{dN}$  elements in  $\mathbb{Z}_q$ , and
- per-query server computation  $2dN$  operations in  $\mathbb{Z}_q$ .

**Fetching many database records at once (“Batch PIR”).** In many applications [9, 10], a client wants to fetch  $k$  records from the PIR server at once. If the client runs our PIR protocol

$k$  times on a database of  $N$  records, the total server time would be roughly  $kN$ . We can apply the “batch PIR” techniques of Ishai et al. [57] to allow a client to fetch  $k$  records at server-side cost  $\ll kN$ , without increasing the hint size.

The idea is to randomly partition the database of  $N$  records into  $k$  chunks, each represented as a matrix of dimension  $(\sqrt{N}/k)$ -by- $\sqrt{N}$ . If the  $k$  records that the client wants to fetch fall into distinct chunks, the client can recover these records by running SimplePIR once on each database chunk. In this case, the hint size remains  $n\sqrt{N}$ —as in one-query SimplePIR. The communication cost for the client is  $k\sqrt{N}$ — $k$  times larger than in one-query SimplePIR (and identical to the communication if the client fetched all  $k$  records sequentially). The server performs  $N$  operations in  $\mathbb{Z}_q$ —as in one-query SimplePIR.

However, more than one of the client’s desired records may fall into the same chunk. There are two ways to handle this:

- If the client must recover all  $k$  records with overwhelming probability, the client can make  $\lambda$  PIR queries to each of the  $k$  chunks to achieve failure probability  $2^{-\Omega(\lambda)}$  [9, 57]. This optimization saves on server work as long as  $\lambda < k$ .
- If the client only needs to recover a constant fraction of the  $k$  database records, then the client and the server can run this batch-PIR protocol only once. The server-side computation cost is as in one-query SimplePIR.

**Additional improvements.** We discuss how to further improve the asymptotic efficiency of SimplePIR in Appendix C.3.

## 4.4 Fast linearly homomorphic encryption

In Appendix D, we introduce the notion of *linearly homomorphic encryption with preprocessing*. This new primitive abstracts out the key properties of Regev encryption that we use in SimplePIR. We expect this new form of linearly homomorphic encryption to have further practical applications.

# 5 DoublePIR

While SimplePIR has high server-side throughput, it requires the client to download and store a relatively large preprocessed hint, of size roughly  $n\sqrt{N}$  on lattice dimension  $n \approx 2^{10}$  and database size  $N$ . In this section, we present DoublePIR, a new PIR scheme that recursively applies SimplePIR to reduce the hint size to roughly  $n^2$  on lattice dimension  $n$ —independent of the database size—while maintaining a server-side throughput upwards of 7.4 GB/s. (In practice, this hint size is 16 MB for one-byte records.) For databases of very many records ( $N \gg n^2 \approx 2^{20}$ ), DoublePIR has a much smaller hint size than SimplePIR. As in SimplePIR, the per-query communication cost for DoublePIR is  $O(\sqrt{N})$  on database size  $N$ .

## 5.1 Construction

We present a formal description of DoublePIR in Figure 14 of Appendix E, along with a full correctness and security analysis.

In this section, we describe the key design ideas.

We first give the concrete costs of DoublePIR on database size  $N$ , lattice dimension  $n$ , LWE modulus  $q$ , plaintext modulus  $p$ , and  $\kappa = \lceil \log(q)/\log(p) \rceil \approx 4$  (chosen as in Appendix E.1). In a one-time public preprocessing phase, DoublePIR requires

1. the server to perform  $2nN + 2\kappa n^2\sqrt{N}$  operations in  $\mathbb{Z}_q$ , and
2. the client to download  $\kappa n^2$  elements in  $\mathbb{Z}_q$ .

On each query, DoublePIR requires

1. the client to upload  $2\sqrt{N}$  elements in  $\mathbb{Z}_q$ ,
2. the server to do  $2N + 2(2n + 1) \cdot \sqrt{N} \cdot \kappa \mathbb{Z}_q$  operations, and
3. the client to download  $(2n + 1) \cdot \kappa$  elements in  $\mathbb{Z}_q$ .

At a high level, DoublePIR first executes exactly as SimplePIR: from the database, the server computes a hint matrix and, in response to each client's query, produces an answer vector. At this point, we observe that a client querying for element  $(i_{\text{row}}, i_{\text{col}})$  in SimplePIR needs two pieces of information to recover its desired database element:

- row  $i_{\text{row}}$  of the hint matrix  $\mathbf{D} \cdot \mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n}$ , and
- element  $i_{\text{row}}$  of the answer vector  $\mathbf{a} \in \mathbb{Z}_q^{\sqrt{N}}$ .

Thus, in DoublePIR, we have the client execute a second level of SimplePIR over the hint matrix and the answer vector to retrieve these  $(n + 1)$  values. As such, the client in DoublePIR recovers the database entry at  $(i_{\text{row}}, i_{\text{col}})$  without downloading the large first-level hint.

Kushilevitz and Ostrovsky [62] first proposed using recursion to reduce communication costs in PIR in this way. However, applied naïvely, this strategy requires  $(n + 1) \approx 2^{10}$  instances of PIR to recover the  $(n + 1)$  desired values. We avoid this bottleneck with the insight that SimplePIR lets the client retrieve a *column* of the database at a time (as discussed in Section 4.3). Therefore, in DoublePIR, we run the second level of PIR over the database corresponding to the *transpose* of the hint matrix concatenated with the answer vector (i.e.,  $[\mathbf{D} \cdot \mathbf{A} \parallel \mathbf{a}]^T$ ). Using a single invocation of SimplePIR, the client in DoublePIR can retrieve column  $i_{\text{row}}$  of this database—which holds exactly row  $i_{\text{row}}$  of the hint matrix and element  $i_{\text{row}}$  of the answer vector—and finally recover the database entry at  $(i_{\text{row}}, i_{\text{col}})$ . As SimplePIR executes over a database of elements in  $\mathbb{Z}_p$ , while the hint matrix and the answer vector consist of elements in  $\mathbb{Z}_q$ , the server in DoublePIR computes the base- $p$  decomposition of the entries in the hint matrix and the answer vector before performing the second level of PIR.

Since this second level of PIR operates on a much smaller database, its cost is dwarfed by that of the first level of PIR: in DoublePIR, both the online communication and the server throughput remain roughly the same as in SimplePIR. Moreover, as the client in DoublePIR forgoes downloading the large first-level hint, it now only downloads a much smaller hint, whose size is independent of the database length, produced by the second level of PIR. Concretely, our PIR client downloads a 16 MB hint in the offline phase.

*Remark 5.1* (Why not recurse more?). DoublePIR performs

two levels of PIR to reduce the total communication. A natural question is whether additional levels of recursion can help, as in standard single-server PIR schemes [62]. After  $r$  levels of recursion, the cost of the recursive PIR scheme, on lattice dimension  $n$  and database size  $N$ , would be (hiding constants):

- one-time download  $n^r$  in the preprocessing step, as well as
- per-query upload  $r \cdot N^{1/r}$  and download  $n^{r-1}$ .

For  $r > 2$ , the communication is likely too large for databases of interest. An intriguing open question is to construct recursive LWE-based PIR schemes with total communication  $n \cdot N^{1/r}$ .

## 5.2 Extensions

We extend DoublePIR to handle diverse deployment scenarios.

**Handling large database records.** To handle databases with large records, we represent each record as a series of elements in  $\mathbb{Z}_p$ , where  $p$  is the plaintext modulus, using base- $p$  decomposition. Let  $d$  denote the number of  $\mathbb{Z}_p$  elements that each record maps to. Then, on each execution of DoublePIR, we run the PIR scheme  $d$  times in parallel, over  $d$  databases, where the  $i$ -th database holds the  $i$ -th  $\mathbb{Z}_p$  element of each record. With this approach, DoublePIR's throughput is identical on databases with long records and with short records. On a database of  $N$  records, each in  $\mathbb{Z}_p^d$ , with lattice dimension  $n$ , LWE modulus  $q$ , and  $\kappa = \lceil \log(q)/\log(p) \rceil$ , DoublePIR has:

- hint size  $d\kappa n^2$  elements in  $\mathbb{Z}_q$ ,
- online upload  $2\sqrt{N}$  elements in  $\mathbb{Z}_q$ ,
- online server work  $2d(N + \kappa(2n + 1)\sqrt{N})$  ops. in  $\mathbb{Z}_q$ , and
- online download  $d\kappa \cdot (2n + 1)$  elements in  $\mathbb{Z}_q$ .

**Batching client queries.** To implement query batching in DoublePIR, we batch queries exactly as in SimplePIR when performing the first level of PIR. As DoublePIR makes non-black-box use of SimplePIR in performing the second level of PIR, we are not able to derive any computation savings from batching in this second, recursive step. (In particular, in the second level of PIR, the client must read an entire column consisting of  $(n + 1)$  elements at once for each query; this breaks SimplePIR's batching trick.) However, as the first level of PIR dominates the computation in DoublePIR, batching many queries nevertheless greatly improves DoublePIR's throughput.

Concretely, to fetch a constant fraction among a set of  $k$  records from a database of  $N$  values in  $\mathbb{Z}_p$ , on lattice dimension  $n$ , LWE modulus  $q$ , and  $\kappa = \lceil \log(q)/\log(p) \rceil$ , DoublePIR has:

- hint size  $\kappa n^2$  elements in  $\mathbb{Z}_q$ ,
- online upload  $\sqrt{N}(k + \sqrt{k})$  elements in  $\mathbb{Z}_q$ ,
- online server work  $2N + 2k(2n + 1)\kappa\sqrt{N}$  ops. in  $\mathbb{Z}_q$ , and
- online download  $k\kappa(2n + 1)$  elements in  $\mathbb{Z}_q$ .



## 6 Data structure for private approximate set membership

In this section, we introduce a new data structure for the *private approximate set membership* problem. In this problem, a client holds a private string  $\sigma$ , a server holds a set of strings  $S$ , and the client wants to test whether  $\sigma \in S$  without revealing  $\sigma$  to the server. Unlike in private set intersection [46], the server’s set  $S$  is public. To rule out the trivial solution where the server sends  $S$  to the client, we insist on communication sublinear in  $|S|$ . Our approach is approximate: there is some chance that the client outputs “ $\sigma \in S$ ” when in fact this is not the case. However, this false-positive rate is bounded even when the set  $S$  and the string  $\sigma$  are chosen *adversarially*. Looking ahead, our data structure will be at the core of our new scheme for auditing in Certificate Transparency (Section 7).

At a high level, we have the server preprocess its set  $S$  into a data structure. Then, the client, holding a string  $\sigma$ , can test whether  $\sigma \in S$  by privately reading a few bits of the server’s data structure using PIR. The relevant cost metrics are:

- *Number of probes.* How many bits of the server’s data structure must the client read?
- *PIR database size.* Over how many bits of the server’s data structure does the client perform its private PIR read?
- *Adversarial false-positive rate.* Given an honest server but an adversarially chosen set  $S$  and string  $\sigma$ , what is the probability, *only over the client’s secret randomness*, that the client outputs “ $\sigma \in S$ ” when in fact  $\sigma \notin S$ ?

*Background: Bloom filters.* A Bloom filter [15] is a standard data structure for approximate set membership. A one-hash-function Bloom filter consists of a fixed-length bitstring  $D$  and uses a hash function  $H: \{0, 1\}^* \rightarrow \{1, \dots, |D|\}$ . Given a set of strings  $S \subseteq \{0, 1\}^*$ , the setup routine hashes each string  $\sigma \in S$  into an index  $i \in \{1, \dots, |D|\}$  and sets the corresponding bit of the data array:  $D_{H(\sigma)} \leftarrow 1$ . Then, to test whether a string  $\sigma$  is in the set represented by the data structure  $D$ , the query algorithm outputs “ $\sigma \in S$ ” if and only if the bit  $D_{H(\sigma)} = 1$ .

As long as the query string is chosen *independently of the hash function  $H$* , the probability of a false-positive is at most  $1/2$  when  $|D| \geq 2|S|$ . However, when the query string is chosen *adversarially*—as can be the case in our application—an adversary can easily find strings  $\sigma \in S$  and  $\hat{\sigma} \notin S$  such that  $H(\sigma) = H(\hat{\sigma})$ . In this case, the one-hash-function Bloom filter will *always* incorrectly output “ $\hat{\sigma} \in S$ .” We present a new data structure—which is a twist on Bloom filters—to address this false-positive issue, without increasing the number of probes or the size of the PIR read required by the query algorithm.

*Remark 6.1* (False positives). Some, but not all, applications can tolerate a non-negligible false-positive rate. For example, credential-breach lookups [66, 82, 91] and contact discovery [60] may tolerate false-positive rates as large as  $2^{-30}$ ; in contrast, Safe Browsing blocklist checks [52,61] demand a cryp-

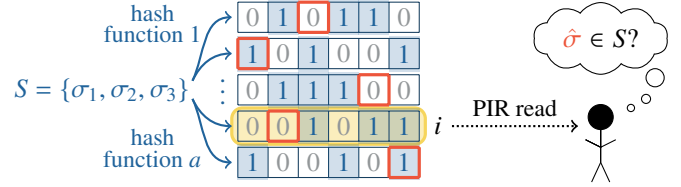


Figure 4: Our data structure for private, approximate set membership with adversarial soundness, when instantiated with a set  $S$  consisting of three strings and with  $a = 5$  hash functions. We highlight in blue the bits of the data structure that are set, in red the bits that the query string  $\hat{\sigma}$  maps to, and in yellow the area covered by the client’s PIR read, when the client probes the  $i$ -th one-hash-function Bloom filter.

topographically negligible false-positive rate, as false positives would cause a legitimate website to be flagged as malicious. In the latter case, other data structures may be more appropriate.

### 6.1 Our approximate membership test

Our data structure for approximate set-membership, illustrated in Figure 4, is parameterized by integers  $a, k \in \mathbb{N}$ , a universe of strings  $\mathcal{U}$ , and a set size  $N$ . The data structure consists of  $a$  independent one-hash-function Bloom filters [15], each of size  $kN$  bits. Crucially, these Bloom filters each use independent hash functions, which are chosen *after* the set  $S$  is fixed. In the remainder of this section, we give an informal description of our construction; a formal treatment appears in Appendix F.

*Data-structure setup.* The setup algorithm takes as input a set of strings  $S \subseteq \mathcal{U}$  of size at most  $N$ . The algorithm then chooses a set of  $a$  hash functions—one per Bloom filter—and inserts each string in  $S$  into each of the  $a$  one-hash-function Bloom filters defined by these hash functions. (In practice, we would use a salted hash function with a different salt per filter.)

*Query algorithm.* Given a query string  $\sigma$ , the query algorithm chooses an index  $i \leftarrow^{\mathcal{R}} [a]$  at random, and outputs the result of querying the  $i$ -th one-hash-function Bloom filter on string  $\sigma$ .

Our data structure has the following properties:

**Correctness.** For any set  $S \subseteq \mathcal{U}$  and any query string  $\sigma \in S$ , the query algorithm always returns “ $\sigma \in S$ .”

**Adversarial false-positive rate 1/2.** For any set  $S \subseteq \mathcal{U}$  of size at most  $N$ , for a random choice of the hash functions used in the data structure, and for any query string  $\hat{\sigma} \notin S$ —which can depend on the hash functions—the data structure incorrectly returns “ $\hat{\sigma} \in S$ ” with probability at most  $1/2$  (taken over the query algorithm’s randomness), for an appropriate choice of the parameters  $a$  and  $k$ . In Appendix F, we prove:

*Proposition 6.2:* For all  $\lambda \in \mathbb{N}$ , on parameters  $k \geq 8$  and  $a \geq 2(\log(|\mathcal{U}|) + \lambda)$ , our approximate set-membership data structure has adversarial false-positive rate at most  $1/2$ . The construction fails with probability  $2^{-\lambda}$ , over the choice of the Bloom filters’ hash functions, modeled as independent random oracles. Concretely, on  $|\mathcal{U}| = 2^{256}$ , taking  $a = 768$  and  $k = 8$  gives false-positive rate  $1/2$  and failure probability  $2^{-128}$ .

**PIR compatibility.** To *privately* test whether a string is in the set, the client can perform a PIR read over only a small fraction of the data structure. More specifically, the query algorithm probes a single bit in one of the Bloom filters. The client can reveal which Bloom filter it wants to probe to the server, as this *does not depend on the query string*. So, while the entire data structure consists of  $akN$  bits, the client can execute a private set-membership test with a PIR read over only  $kN$  bits.

## 6.2 Related approaches and comparison

We now compare our solution to other data structures for private approximate set-membership, given in Table 5. One natural alternative would be to use a *single* one-hash-function Bloom filter. (In contrast, our construction uses  $a \approx 768$  one-hash-function Bloom filters.) However, this approach is *not* sound in our adversarial setting: as the data structure (including its hash function) is public, an adversary can trivially find a string that causes the query algorithm to *always* return a false-positive result. We can address this issue by using a Bloom filter with  $O(\lambda)$  hash functions, which gives security against  $2^{O(\lambda)}$ -time attacks (where  $\lambda \approx 128$  is a security parameter). Unfortunately, the query algorithm of such a data structure is roughly  $\lambda \times$  more expensive than ours in terms of both (1) the number of probes and (2) the size of the PIR read required for a private query.

Adversarial Bloom filters [28, 45, 78] provide the false-positive guarantees we require, but do not naturally support private reads via PIR. In particular, they have the client send its string  $\sigma$  to the server; the server then applies a pseudorandom function to  $\sigma$  to determine which bits to probe. It is not clear how to use such a data structure in our setting without relatively expensive general-purpose multi-party computation schemes.

Another approach to private set membership has the client and the server execute a PIR by keywords protocol [26]. On security parameter  $\lambda$ , the server stores a  $\lambda$ -bit hash of each string in its set in a hash table. Thereafter, the client can perform PIR over this hash table to check if a string is present. While the client probes only few locations of the hash table, its PIR read must cover the entire table, or roughly  $3\lambda N$  bits.

Finally, prior work constructs other data structures for approximate set membership [36, 53], offering different performance trade-offs. Combining our ideas for efficiently tolerating false positives in an adversarial setting with such data structures is an intriguing direction for future work.

## 7 Application: Auditing in Certificate Transparency

We now apply our new PIR schemes (Sections 4 and 5), along with our set-membership data structure (Section 6), to solve the problem of privately auditing signed certificate timestamps in deployments of Certificate Transparency [63, 64, 74].

	Probes	PIR size	Adv. false-positive rate
PIR by keywords [26]	2	$3\lambda N$	0
Standard Bloom filter [15]	$O(\lambda)$	$O(\lambda N)$	0
1-hash-fn Bloom filter [15]	1	$2N$	1 (insecure)
<b>This work</b>	1	$8N$	1/2

Table 5: Private set-membership data structures, for sets of  $N$  elements from universe  $\mathcal{U}$ , on security parameter  $\lambda$ . The data structures may fail with probability  $|\mathcal{U}| 2^{-\lambda}$ , over their random choice of hash functions, modeled as random oracles.

### 7.1 Problem statement

**Background: Certificate Transparency.** The goal of Certificate Transparency is to store every public-key certificate that every certificate authority issues in a set of publicly accessible logs. To this end, certificate authorities submit the certificates they issue to log operators, who respond with a *signed certificate timestamp* (SCT). The SCT is a promise to include the new certificate in the log maintained by this operator within some bounded period of time.

Later on, when a TLS server sends a public-key certificate to a client, the server attaches a number of SCTs according to the client’s policy (e.g., Chrome and Safari both require SCTs from three distinct log operators). By verifying the SCTs, the client can be sure that each of the log operators has seen the new certificate and—if the operator is honest—will eventually log it. Domain operators can then use the logs to detect whether a certificate authority has mistakenly or maliciously issued a certificate for their domain. In this setting, the log contents are public; related work investigates scenarios where this is not the case, as in end-user key distribution [75].

**SCT auditing.** To keep the logs honest, some party in the system must verify that the log operators are fulfilling the promise implicit in the SCTs that they issue. In particular, if a client receives an SCT for some certificate  $C$  signed by a log operator, the client would like to verify that  $C$  appears in that operator’s public log. This process is *SCT auditing*.

Clients must be involved in SCT auditing, as they are the only participants who see SCTs “in the wild.” However, the set of SCTs that a client sees reveals information about the client’s browsing history: the fact that a client has seen an SCT for `example.com` reveals that the client has visited `example.com`. Thus, to protect its privacy, the client should not reveal which SCTs it has seen to the log operators or to any other entity.

Google’s recent solutions for SCT auditing [35, 89] involve an *SCT auditor* (run by Google) that is separate from the client. In their model, the auditor maintains the entire set of SCTs for non-expired certificates from all Certificate Transparency logs. Every SCT that a client sees for a live website should appear in the auditor’s set. To determine whether an SCT is valid, a client can check whether it (or really, its SHA256 hash) appears in the set of valid SCTs maintained by the auditor:

- If the client’s SCT appears in the auditor’s set, then the log

server that issued the SCT correctly fulfilled its promise.

- If not, the client can report the problematic SCT to the auditor to investigate further. Prior work shows how this can be done while keeping the SCT in question hidden [41].

A privacy-protecting solution for SCT auditing must allow the client to test whether its SCT appears in the auditor’s set, *without* revealing its SCT to the auditor. This is a private set-membership problem [91]. While on its surface this problem resembles other applications of PIR in the literature [59, 71], the fact that many clients engage in the protocol with the same auditor means that we can tolerate false positives. That is, it is acceptable for a client to incorrectly believe that an SCT is in the auditor’s set, since over many clients we can expect that missing SCTs are eventually identified. To summarize, we require the following properties, which we state only informally:

- **Correctness with false positives.** When an honest client holding string  $\sigma$ , chosen independently of the client’s secret randomness, interacts with an honest auditor holding set  $S$ :
  - if  $\sigma \in S$ , then the client always outputs “valid,” and
  - if  $\sigma \notin S$ , then the client outputs “valid” with probability at most  $1/2$ , over the choice of the client’s randomness.
- **Privacy for the client.** When an honest client interacts with a malicious auditor, the auditor learns nothing about the client’s private input string  $\sigma$ .

We do not require correctness to hold against a malicious auditor; such an auditor could trivially lie about its set of SCTs.

**System parameters.** There are roughly five billion active SCTs in the web today [35]. Roughly six million of these are added or removed each day as certificate authorities issue certificates and as certificates expire [3]. Google Chrome’s current proposal for SCT auditing has a *false-positive rate* of essentially zero: when a client audits an SCT, it correctly learns whether the SCT is valid. However, Chrome’s proposal has a *detection rate* of  $1/10,000$ : the Chrome client randomly samples  $0.01\%$  of the SCTs associated with its TLS connections, and audits only this small fraction of all SCTs [2]. This random sampling reduces the amortized cost of auditing by  $10,000\times$ , but also reduces the chance that any single auditing client catches a cheating log. Still, across many auditing clients, this randomized SCT auditing catches—with high probability—widely distributed invalid SCTs. After  $10,000$  clients observe an invalid SCT, in expectation one will audit it and implicate the cheating log.

**Existing approaches.** Two notable proposals for SCT auditing—which do not provide cryptographic privacy—are:

*Opt-out SCT auditing.* Chrome’s current approach [35] has the client reveal the first 20 bits of the hash of its SCT to the auditor [1]. The auditor replies with all roughly 1000 SCTs in its set that match this 20-bit prefix. This method achieves  $k$ -anonymity for  $k = 1000$ , i.e., it leaks that the client visited one of a set of 1000 sites.

*Anonymizing proxy.* The client could use proxy servers, such as in Tor [37], to send its SCT to the auditor anonymously [32];

the auditor could reply with the bit indicating whether the SCT appears in its set. This mechanism reveals the entire distribution of clients’ SCTs to the auditor and is susceptible to timing attacks, which could allow the auditor to deanonymize particular clients.

## 7.2 Our approach

We propose a new scheme for SCT auditing that achieves cryptographic privacy. The deployment is as follows:

*1. Auditor: Data-set construction.* The auditor prepares an approximate set-membership data structure holding all SHA256 hashes of all  $N$  active SCTs. This data structure consists of  $a = 768$  arrays, each  $8N$  bits in length, and has false-positive rate  $\epsilon = 1/2$  (Proposition F.1). Then, the auditor runs the PIR Setup routine on each of these  $a$  arrays, producing  $a$  PIR hints.

*2. Client: Hint download.* The client chooses a secret, random index  $i^* \leftarrow [a]$  and downloads the  $i^*$ -th hint from the auditor, revealing  $i^*$  to the auditor in the process. Whenever the client wants to test whether some SCT appears in the auditor’s set, the client can now read a single bit from the auditor’s  $i^*$ -th array. The probability that a cheating log can trick the client into accepting an invalid SCT is at most  $\epsilon$ , the false-positive rate of the underlying set-membership data structure.

If the client audits an  $f$ -fraction of all of its TLS connections, the detection rate is  $f \cdot (1 - \epsilon)$ . In our deployment, we take  $f = 1/5,000$  and  $\epsilon = 1/2$ . This choice gives an overall detection rate of  $1/10,000$ —matching that of Chrome’s current approach [35].

*3. Client and auditor: SCT lookup via PIR.* Each time the client decides to audit an SCT, the client computes the bit of the auditor’s  $i^*$ -th array that it needs to check to verify the SCT’s validity. The client reads this bit privately by running the PIR protocol’s online phase with the auditor, over the  $i^*$ -th array.

In this approach, the client reuses the same secret index  $i^*$  for multiple SCT lookups. As a result, the events that a client fails to detect invalid SCT  $A$  and invalid SCT  $B$  are correlated, whereas in Chrome today these events are independent. However, in our approach, the probability that a client catches the *first* invalid SCT that it looks up remains  $1 - \epsilon = 1/2$ . As such, the client will catch *at least one* invalid SCT with probability  $1/10,000$  and will thus implicate a cheating log with these odds. Any log that cheats more than  $10,000$  clients will be caught in expectation.

**Database updates.** In our proposal, the client holds a PIR hint that depends on the auditor’s set of active SCTs. Whenever this set changes, which happens continuously as certificates are issued and expire, the auditor must update its set-membership data structure. Without extra engineering, the client would have to download a fresh PIR hint from the auditor each time.

Our approach has the client download a fresh hint only periodically—once per month, for example. The auditor specifies the range of *certificate issue dates* that each hint covers. When the client decides to audit an SCT, it checks whether its current hint covers the issue date of the SCT in question. If so,

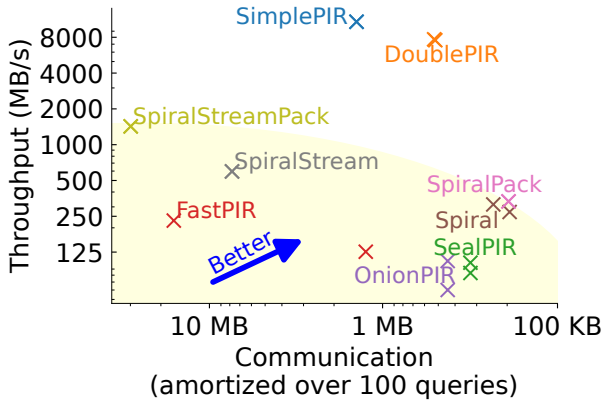


Figure 6: Throughput vs. per-query communication, on a 1 GB database. For each PIR scheme, we display the communication and the corresponding throughput for two choices of entry size: one that maximizes throughput, and another that minimizes communication. (For schemes displayed only once, both entry sizes are the same.) The communication cost is the total (i.e., offline and online) communication, amortized over 100 queries. We highlight prior work in yellow.

the client tests the SCT’s validity; if not, the client caches the SCT so as to test its validity the next time it downloads a hint. In this way, the client eventually audits its full random sample of SCTs, but reuses each hint for multiple SCT lookups. The server must now store multiple versions of the database, which is relatively inexpensive; in a large-scale deployment, one or more physical servers could hold each version in memory.

## 8 Evaluation

**Implementation.** We implement SimplePIR in fewer than 1,200 lines of Go code, along with 200 lines of C, and DoublePIR in 210 additional lines of Go code. Our code does not rely on any external libraries and is published under the MIT open-source license at [github.com/ahenzinger/simplepir](https://github.com/ahenzinger/simplepir). We give sample code for SimplePIR in Appendix G.

We use the appropriate data types to natively support operations over  $\mathbb{Z}_q$  (e.g., `uint32` for  $q = 2^{32}$ ). We store the database in memory in packed form and decompress it into  $\mathbb{Z}_p$  elements on-the-fly, as otherwise the Answer routine is memory-bandwidth-bound. In DoublePIR, we represent the database as a rectangular (rather than square) matrix, so that the first level of PIR dominates the computation.

We run all experiments using a single thread of execution, on an AWS `c5n.metal` instance running Ubuntu 22.04. To collect the throughput numbers for tables, we run each scheme five times and report the average. All standard deviations in throughput are smaller than 10% of the throughput measured.

	Communication				Throughput (MB/s)
	Offline (MB)		Online (KB)		
	Up.*	Down.*	Up.	Down.	
SealPIR	5	0	91	181	97
FastPIR	0.06	0	33 000	64	217
OnionPIR	5	0	256	128	60
Spiral	15	0	14	20	259
SpiralPack	19	0	14	20	260
SpiralStream	0.34	0	15 000	20	485
SpiralStreamPack	15	0	29 000	99	1,370*
<b>SimplePIR</b> (\$4)	0	121	121	121	10,138
<b>DoublePIR</b> (\$5)	0	16	313	32	7,622

Table 8: PIR scheme performance on a database of  $2^{33} \times 1$ -bit entries. We highlight in green cells that are within  $5\times$  of the best, and in red cells that are within  $5\times$  of the worst, in their respective columns. (We leave uncolored cells that are within  $5\times$  of the best and worst.) We automatically “re-balance” schemes without an automatic parameter selection tool (SealPIR, FastPIR, and OnionPIR), by executing them on a database of  $2^{33}/d$  entries, each of size  $d$ , where  $d$  is the closest valid power-of-2 to the scheme’s “optimal” entry size (see Table 10). \*The offline upload is equal to the per-client server storage. \*The offline download is equal to the client storage. \*The throughput here is slightly higher than in Table 1 due to variance in the measurements.

### 8.1 Microbenchmarks

**Throughput.** We first measure the maximal throughput of each PIR scheme, on the database dimensions that suit it best. In Table 1, we report the throughput measured for each scheme, on a database roughly 1 GB in size, where we take the entry size to be that for which the highest throughput was reported in the corresponding paper (or in a related paper, if it is not made explicit). These entry sizes appear in Table 10. We confirm that these throughputs are indeed the best achievable by measuring each scheme’s throughput on each entry size in Figure 7.

SimplePIR and DoublePIR achieve throughputs of 10.0 GB/s and 7.4 GB/s respectively, which is roughly  $8\times$  faster than the best prior single-server PIR scheme designed for the streaming setting (SpiralStreamPack) and  $30\times$  faster than the best prior single-server PIR scheme designed for databases with short entries (Spiral). SimplePIR and DoublePIR exceed the *per-server* throughput<sup>1</sup> of some prior two-server PIR schemes: two-server PIR from DPFs [59] has a throughput of 5.3 GB/s/core. Finally, we benchmark the throughput of performing only XORs over a database to provide a hard upper bound on the speed of linear-work, two-server PIR [13,27]. When each server performs a linear scan of XORs over the database, two-server PIR’s throughput is 5.9 GB/s/core. When each server performs a linear scan of XORs over a *random half* of the database, two-server PIR’s throughput is 11.5 GB/s/core—but this requires a non-constant-time implementation (see discussion in Table 1).

**Communication.** In Figure 7, we give each scheme’s total

<sup>1</sup>In computing the *per-server* throughput of two-server PIR (from DPFs and from XOR), we divide the measured throughput by two.

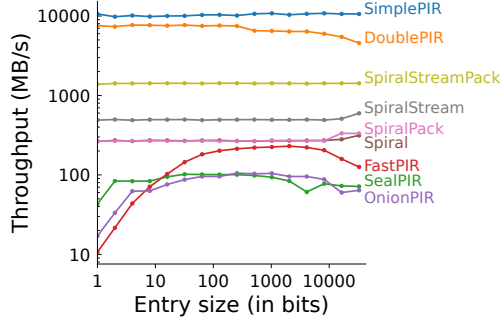
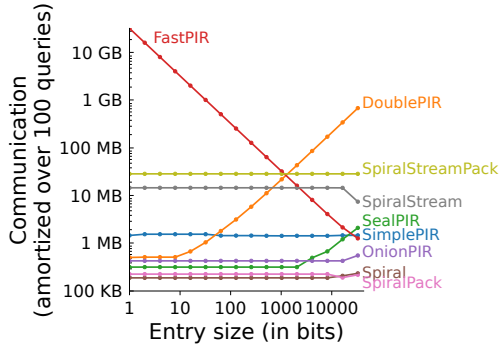


Figure 7: Throughput and per-query communication for each PIR scheme, on a 1 GB database with entries of increasing size. The per-query communication cost is the total (i.e., offline and online) per-query communication, amortized over 100 queries.

communication, amortized over 100 queries, for increasing entry sizes. On databases with short entries, DoublePIR’s amortized communication is comparable to that of the most communication-efficient schemes (Spiral, SpiralPack, SealPIR, and OnionPIR). With larger entries, DoublePIR’s amortized communication costs increase, as the client must download many hints. The two schemes with the closest throughput to ours (SpiralStream and SpiralStreamPack), as well as FastPIR, have much larger amortized communication than both DoublePIR and SimplePIR on entry sizes less than a kilobit.

**Throughput vs. communication trade-off.** We summarize these findings in Figure 6, which displays the throughput/communication trade-off achieved by each PIR scheme. Concretely, we run each scheme on a database of  $2^{33}$  bits with increasing entry sizes (as also done in Figure 7). Then, for each scheme, we display the per-query communication (amortized over 100 queries) and the corresponding throughput for two choices of the entry size: one that maximizes the throughput, and another that minimizes the communication. Figure 6 demonstrates that our new PIR schemes achieve a novel point in the design space: SimplePIR and DoublePIR have substantially higher throughput than all prior single-server PIR schemes; DoublePIR further has a per-query communication cost that is competitive with the most communication-efficient schemes.

**Comparison on a database of  $2^{33} \times 1$ -bit entries.** In Table 8, we give a fine-grained comparison of the performance of each scheme on a database relevant to our application, consisting of  $2^{33}$  1-bit entries. On this database, SimplePIR and DoublePIR again achieve much higher throughput than all other schemes (9.9 GB/s and 7.4 GB/s respectively). SimplePIR has high offline download and thus also client-side storage costs. However, DoublePIR’s offline download is comparable to the offline communication of other PIR schemes, and its online communication is on the order of kilobytes.

For each scheme, we additionally compute its cost per query, when the client makes 100 queries, using the AWS costs for compute ( $\$1.5 \cdot 10^{-5}$ /core-second) and data transfer out of Amazon EC2 ( $\$0.09$ /GB). SimplePIR’s per-query cost is  $\$1 \cdot 10^{-4}$ , while DoublePIR and the cheapest scheme from related work (SpiralStreamPack) each achieve a per-query cost

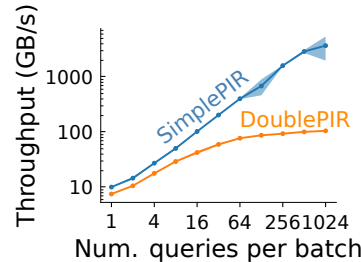


Figure 9: Effective PIR throughput (database size  $\times$  queries per second), with increasing batch sizes and a fixed-size hint, on a database consisting of  $2^{33} \times 1$ -bit entries. The shading displays the standard deviation.

of  $\$2 \cdot 10^{-5}$ . We note, however, that SpiralStreamPack requires megabytes of online upload, which is not reflected in its per-query cost, as AWS only charges for outgoing communication.

**Batching queries.** Finally, we evaluate how SimplePIR and DoublePIR’s effective throughput scales when the client makes a batch of queries for  $k$  records at once, assuming the client only needs to recover a constant fraction of the  $k$  records. For increasing values of  $k$ , we compute the expected number of “successful” queries (i.e., the expected number of queries that fall into a distinct database chunk, as discussed in Section 4.3) and we derive the expected “successful” throughput—that is, the throughput measured when the server answers that number of queries at once, with a single pass over the database.

Figure 9 shows that SimplePIR and DoublePIR’s throughput increases when the client makes a batch of queries at once. SimplePIR’s throughput scales linearly, achieving a value of over 100 GB/s on batch size  $k \geq 16$  and 1000 GB/s on batch size  $k \geq 256$ . DoublePIR achieves a throughput over 50 GB/s for  $k \geq 32$ ; when  $k \geq 256$ , the throughput plateaus at roughly 100 GB/s, as the second level of PIR becomes a bottleneck.

In ??, we give additional benchmarks that measure the server preprocessing time, the client time, and the non-amortized communication of our new PIR schemes, along with tables containing the data displayed in Figures 6, 7 and 9.

## 8.2 Certificate Transparency benchmark

We propose using DoublePIR for the SCT auditing application. With our new data structure for private set membership, the task of SCT auditing requires a single round of PIR over a database with 1-bit entries. For such a database, our microbenchmarks

in Section 8.1 show that DoublePIR achieves both high server throughput and small client storage and communication. SCT auditing occurs in the background, and is not on the critical path to web browsing. Thus, while using PIR may increase the latency of auditing, we believe this is a desirable trade-off in exchange for cryptographic privacy, as long as the computation remains modest and the communication remains comparable.

To benchmark DoublePIR in this application context, we evaluate the scheme on a database consisting of  $2^{36} \times 1$ -bit entries, which is the size of a Bloom filter in our approximate set membership data structure when we instantiate it with all 5 billion active SCTs. (In our evaluation, each entry is a random bit.) On this database size, we measure that DoublePIR has a “hint” of size 16 MB, an online upload of 724 KB, and an online download of 32 KB. The server can answer each query in fewer than 1.3 core-seconds (and this work is fully parallelizable).

As our client must audit one in every 5,000 TLS connections, our proposal for SCT auditing then requires: (1) 16 MB of client storage and download every month (to keep the client hint), and (2) per TLS connection, an amortized overhead of 0.0003 core-seconds of server compute and 150 bytes of communication. Using the AWS costs for compute ( $\$1.5 \cdot 10^{-5}$  per core second) and data transfer ( $\$0.09$  per outgoing GB), for each client, this amounts to a fixed cost of  $\$0.001$  per month, along with  $\$4 \cdot 10^{-9}$  per TLS connection. For a typical client making  $10^4$  TLS connections per week [4], we expect this cost to be roughly  $\$0.02$  per year. Since (after sampling its TLS connections) a typical client makes only few queries to the auditor using each month’s hint, we can reduce the client’s storage to less than 150 KB using the optimization from Appendix E.3.

By comparison, Chrome’s SCT auditing scheme [35] provides only  $k$ -anonymity for  $k = 1,000$ : the server learns that a client visited one of a set of 1,000 domains. Auditing incurs an amortized overhead of 24 B of communication per connection<sup>2</sup>, negligible server computation, and no client storage (unless the client caches popular SCTs). Again assuming a client making  $10^4$  TLS connections per week [4], we expect this scheme to cost roughly  $\$0.001$ /client/year. Our approach using DoublePIR incurs 23× more communication to achieve the goal of cryptographic privacy.

## 9 Conclusion

We show that the per-core throughput of single-server PIR can approach the memory bandwidth of the machine and the performance of two-server PIR. Two exciting directions remain open: one is to reduce our schemes’ communication; another is to combine our ideas with those of *sublinear*-time PIR [30, 31] to reduce the computation beyond the linear-server-time barrier.

<sup>2</sup>The communication may be smaller in practice as the total number of SCTs and the number of SCTs matching each prefix varies.

**Acknowledgements.** We thank Martin Albrecht for answering questions about LWE hardness estimates, Vadim Lyubashevsky for advice on discrete gaussian sampling, and Adam Belay and Zhenyuan Ruan for discussions about AVX performance. We are grateful to Anish Athalye, Derek Leung, and Ryan Lehmkuhl for reviewing a draft of this work, and to Dima Kogan, David Wu, Jean-Philippe Bossuat, Samir Menon, Alex Davidson, Sofía Celi, Emily Stark, Kevin Yeo, and Joe DeBlasio for helpful conversations and feedback. We thank Sebastian Angel for constructive comments on the discussion of malicious security in an earlier version of this work, and for suggestions on how to improve the presentation. We thank Yuval Ishai and Matan Hamilis for discussing how best to compare against two-server PIR schemes. Yuval also helpfully suggested the “XOR PIR fast” construction discussed in Table 1. This work was supported in part by the National Science Foundation (Award CNS-2054869), a gift from Google, a Facebook Research Award, and MIT’s Fintech@CSAIL Initiative. Alexandra Henzinger was supported by the National Science Foundation Graduate Research Fellowship under Grant No. 2141064 and an EECS Great Educators Fellowship. Matthew M. Hong was funded by NIH R01 HG010959. Vinod Vaikuntanathan was supported by DARPA under Agreement No. HR00112020023, NSF CNS-2154149, MIT-IBM Watson AI, Analog Devices, a Microsoft Trustworthy AI grant and a Thornton Family Faculty Research Innovation Fellowship. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the United States Government or DARPA.

## References

- [1] Chrome SCT Auditing Hash Prefix Length. [https://source.chromium.org/chromium/chromium/src/+main:services/network/sct\\_auditing/sct\\_auditing\\_reporter.cc;l=39;drc=ab2806d582b36d2da8d2178aca83031ab99ed5da](https://source.chromium.org/chromium/chromium/src/+main:services/network/sct_auditing/sct_auditing_reporter.cc;l=39;drc=ab2806d582b36d2da8d2178aca83031ab99ed5da). Accessed 4 Nov 2022.
- [2] Chrome SCT Auditing Sampling Rate. [https://source.chromium.org/chromium/chromium/src/+main:chrome/common/chrome\\_features.cc;l=1007;drc=50d8da971873550eb909b9c177cf6188e81ff4c3](https://source.chromium.org/chromium/chromium/src/+main:chrome/common/chrome_features.cc;l=1007;drc=50d8da971873550eb909b9c177cf6188e81ff4c3). Accessed 4 Nov 2022.
- [3] Merkle town. <https://merkle.town/>.
- [4] Mozilla Telemetry Portal, Measurement Dashboard. [https://telemetry.mozilla.org/new-pipeline/dist.html#!cumulative=0&end\\_date=2022-07-17&include\\_spill=0&keys=\\_\\_none\\_\\_!\\_\\_none\\_\\_!\\_\\_none\\_\\_&max\\_channel\\_version=nightly%252F104&measure=HTTP\\_TRANSACTION\\_IS\\_SSL&min\\_channel\\_version=nightly%252F104&processType=\\*&product=Firefox&sanitize=1&sort\\_by=value](https://telemetry.mozilla.org/new-pipeline/dist.html#!cumulative=0&end_date=2022-07-17&include_spill=0&keys=__none__!__none__!__none__&max_channel_version=nightly%252F104&measure=HTTP_TRANSACTION_IS_SSL&min_channel_version=nightly%252F104&processType=*&product=Firefox&sanitize=1&sort_by=value)

- 0&sort\_keys=submissions&start\_date=2022-06-27&table=1&trim=1&use\_submission\_date=0. Accessed 19 July 2022.
- [5] Carlos Aguilar-Melchor, Joris Barrier, Laurent Fousse, and Marc-Olivier Killijian. XPIR: Private information retrieval for everyone. *PoPETs*, 2016.
- [6] Ishtiyaque Ahmad, Yuntian Yang, Divyakant Agrawal, Amr El Abbadi, and Trinabh Gupta. Addra: Metadata-private voice communication over fully untrusted infrastructure. In *OSDI*, 2021.
- [7] Martin Albrecht, Rachel Player, and Sam Scott. On the concrete hardness of learning with errors. In *Journal of Mathematical Cryptology*, 2015.
- [8] Asra Ali, Tancrede Lepoint, Sarvar Patel, Mariana Raykova, Phillipp Schoppmann, Karn Seth, and Kevin Yeo. Communication–Computation trade-offs in PIR. In *USENIX Security*, 2021.
- [9] Sebastian Angel, Hao Chen, Kim Laine, and Srinath Setty. PIR with compressed queries and amortized query processing. In *S&P*, 2018.
- [10] Sebastian Angel and Srinath Setty. Unobservable communication over fully untrusted infrastructure. In *OSDI*, 2016.
- [11] Michael Backes, Aniket Kate, Matteo Maffei, and Kim Pecina. ObliviAd: provably secure and practical online behavioral advertising. In *S&P*, 2012.
- [12] W. Banaszczyk. Inequalities for convex bodies and polar reciprocal lattices in  $\mathbf{R}_n$ . *Discrete & computational geometry*, 1995.
- [13] Amos Beimel, Yuval Ishai, and Tal Malkin. Reducing the servers’ computation in private information retrieval: PIR with preprocessing. *J. Cryptol.*, 2004.
- [14] Mihir Bellare and Phillip Rogaway. Random oracles are practical: A paradigm for designing efficient protocols. In *CCS*, 1993.
- [15] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 1970.
- [16] Elette Boyle, Niv Gilboa, and Yuval Ishai. Function secret sharing: Improvements and extensions. In *CCS*, 2016.
- [17] Elette Boyle, Yuval Ishai, Rafael Pass, and Mary Wootters. Can we access a database both locally and privately? In *TCC*, 2017.
- [18] Elette Boyle, Lisa Kohl, and Peter Scholl. Homomorphic secret sharing from lattices without fhe. In *EUROCRYPT*, 2019.
- [19] Zvika Brakerski. Fully homomorphic encryption without modulus switching from classical GapSVP. In *CRYPTO*, 2012.
- [20] Zvika Brakerski, Nico Döttling, Sanjam Garg, and Giulio Malavolta. Leveraging linear decryption: Rate-1 fully-homomorphic encryption and time-lock puzzles. In *TCC*, 2019.
- [21] Christian Cachin, Silvio Micali, and Markus Stadler. Computationally private information retrieval with poly-logarithmic communication. In *EUROCRYPT*, 1999.
- [22] Ran Canetti, Justin Holmgren, and Silas Richelson. Towards doubly efficient private information retrieval. In *TCC*, 2017.
- [23] Yan-Cheng Chang. Single database private information retrieval with logarithmic communication. In *ACISP*, 2004.
- [24] Melissa Chase, Sanjam Garg, Mohammad Hajiabadi, Jialin Li, and Peihan Miao. Amortizing rate-1 OT and applications to PIR and PSI. In *TCC*, 2021.
- [25] Massimo Chenal and Qiang Tang. On key recovery attacks against existing somewhat homomorphic encryption schemes. Cryptology ePrint Archive, Report 2014/535, 2014.
- [26] Benny Chor, Niv Gilboa, and Moni Naor. Private information retrieval by keywords. Cryptology ePrint Archive, Report 1998/003, 1998.
- [27] Benny Chor, Oded Goldreich, Eyal Kushilevitz, and Madhu Sudan. Private information retrieval. In *FOCS*, 1995.
- [28] David Clayton, Christopher Patton, and Thomas Shrimpton. Probabilistic data structures in adversarial environments. In *CCS*, 2019.
- [29] Don Coppersmith and Shmuel Winograd. Matrix multiplication via arithmetic progressions. In *STOC*, 1987.
- [30] Henry Corrigan-Gibbs, Alexandra Henzinger, and Dmitry Kogan. Single-server private information retrieval with sublinear amortized time. In *EUROCRYPT*, 2022.
- [31] Henry Corrigan-Gibbs and Dmitry Kogan. Private information retrieval with sublinear online time. In *EUROCRYPT*, 2020.

- [32] Rasmus Dahlberg, Tobias Pulls, Tom Ritter, and Paul Syverson. Privacy-preserving and incrementally-deployable support for Certificate Transparency in Tor. *PoPETS*, 2021.
- [33] Ivan Damgård and Mads Jurik. A generalisation, a simplification and some applications of Paillier’s probabilistic public-key system. In *PKC*, 2001.
- [34] Alex Davidson, Gonçalo Pestana, and Sofia Celi. Frodopir: Simple, scalable, single-server private information retrieval. Cryptology ePrint Archive, Paper 2022/981, 2022.
- [35] Joe DeBlasio. Opt-out SCT auditing in Chrome. <https://docs.google.com/document/d/16G-Q7iN3kB46GSW5b-sfH5M03nKSYyEb77Ysm7TMZGE/edit>.
- [36] Peter C. Dillinger and Stefan Walzer. Ribbon filter: practically smaller than bloom and xor. *CoRR*, abs/2103.02515, 2021.
- [37] Roger Dingledine, Nick Mathewson, and Paul Syverson. Tor: The Second-Generation onion router. In *USENIX Security*, 2004.
- [38] Yevgeniy Dodis, Shai Halevi, Ron D. Rothblum, and Daniel Wichs. Spooky encryption and its applications. In *CRYPTO*, 2016.
- [39] Nico Döttling, Sanjam Garg, Yuval Ishai, Giulio Malavolta, Tamer Mour, and Rafail Ostrovsky. Trapdoor hash functions and their applications. In *CRYPTO*, 2019.
- [40] Taher ElGamal. A public key cryptosystem and a signature scheme based on discrete logarithms. *IEEE transactions on information theory*, 1985.
- [41] Saba Eskandarian, Eran Messeri, Joseph Bonneau, and Dan Boneh. Certificate Transparency with privacy. In *PETS*, 2017.
- [42] Junfeng Fan and Frederik Vercauteren. Somewhat practical fully homomorphic encryption. Cryptology ePrint Archive, Paper 2012/144, 2012.
- [43] Prastudy Fauzi, Martha Norberg Hovd, and Håvard Raddum. A practical adaptive key recovery attack on the LGM (GSW-like) cryptosystem. In *PQCrypto*, 2021.
- [44] Prastudy Fauzi, Martha Norberg Hovd, and Håvard Raddum. On the IND-CCA1 security of FHE schemes. *Cryptography*, 2022.
- [45] Mia Filić, Kenneth G. Paterson, Anupama Unnikrishnan, and Fernando Virdia. Adversarial correctness and privacy for probabilistic data structures. *CCS*, 2022.
- [46] Michael J. Freedman, Kobbi Nissim, and Benny Pinkas. Efficient private matching and set intersection. In *EUROCRYPT*, 2004.
- [47] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *STOC*, 2009.
- [48] Craig Gentry and Shai Halevi. Compressible FHE with applications to PIR. In *TCC*, 2019.
- [49] Craig Gentry and Zulfikar Ramzan. Single-database private information retrieval with constant communication rate. In *ICALP*, 2005.
- [50] Niv Gilboa and Yuval Ishai. Distributed point functions and their applications. In *EUROCRYPT*, 2014.
- [51] Shafi Goldwasser and Silvio Micali. Probabilistic encryption. *Journal of Computer and System Sciences*, 1984.
- [52] Google. Safe Browsing APIs (v4). <https://developers.google.com/safe-browsing/v4>.
- [53] Thomas Mueller Graf and Daniel Lemire. Xor filters: Faster and smaller than bloom and cuckoo filters. *ACM J. Exp. Algorithmics*, 2020.
- [54] Matthew Green, Watson Ladd, and Ian Miers. A protocol for privately reporting ad impressions at scale. In *CCS*, 2016.
- [55] Trinabh Gupta, Natacha Crooks, Whitney Mulhern, Srinath Setty, Lorenzo Alvisi, and Michael Walfish. Scalable and private media consumption with Popcorn. In *NSDI*, 2016.
- [56] Daniel Günther, Maurice Heymann, Benny Pinkas, and Thomas Schneider. GPU-accelerated PIR with Client-Independent preprocessing for Large-Scale applications. In *Usenix Security*, 2022.
- [57] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Batch codes and their applications. In *STOC*, 2004.
- [58] Ari Juels. Targeted advertising ... and privacy too. In *CT-RSA*, 2001.
- [59] Daniel Kales, Olamide Omolola, and Sebastian Ramacher. Revisiting user privacy for certificate transparency. In *EuroS&P*, 2019.
- [60] Daniel Kales, Christian Rechberger, Thomas Schneider, Matthias Senker, and Christian Weinert. Mobile private contact discovery at scale. In *USENIX Security*, 2019.
- [61] Dmitry Kogan and Henry Corrigan-Gibbs. Private blocklist lookups with Checklist. In *USENIX Security*, 2021.



- [62] Eyal Kushilevitz and Rafail Ostrovsky. Replication is not needed: Single database, computationally-private information retrieval. In *FOCS*, 1997.
- [63] Ben Laurie. Certificate transparency. *Communications of the ACM*, 2014.
- [64] Ben Laurie, Adam Langley, and Emilia Kasper. Certificate transparency. RFC 6962, 2013.
- [65] François Le Gall. Faster algorithms for rectangular matrix multiplication. In *FOCS*, 2012.
- [66] Lucy Li, Bijeeta Pal, Junade Ali, Nick Sullivan, Rahul Chatterjee, and Thomas Ristenpart. Protocols for checking compromised credentials. In *CCS*, 2019.
- [67] Jilan Lin, Ling Liang, Zheng Qu, Ishtiyaque Ahmad, Liu Liu, Fengbin Tu, Trinabh Gupta, Yufei Ding, and Yuan Xie. INSPIRE: In-storage private information retrieval via protocol and architecture co-design. In *ISCA*, 2022.
- [68] Richard Lindner and Chris Peikert. Better key sizes (and attacks) for LWE-based encryption. In *CT-RSA*, 2011.
- [69] Helger Lipmaa. An oblivious transfer protocol with log-squared communication. In *International Conference on Information Security*, 2005.
- [70] Jake Loftus, Alexander May, Nigel P. Smart, and Frederik Vercauteren. On CCA-secure somewhat homomorphic encryption. In *Selected Areas in Cryptography*, 2012.
- [71] Wouter Lueks and Ian Goldberg. Sublinear scaling for multi-client private information retrieval. In *International Conference on Financial Cryptography and Data Security*, 2015.
- [72] Vadim Lyubashevsky, Chris Peikert, and Oded Regev. On ideal lattices and learning with errors over rings. *Journal of the ACM*, 2013.
- [73] Yiping Ma, Ke Zhong, Tal Rabin, and Sebastian Angel. Incremental offline/online PIR. In *USENIX Security*, 2022.
- [74] Sarah Meiklejohn, Joe DeBlasio, Devon O’Brien, Chris Thompson, Kevin Yeo, and Emily Stark. SoK: SCT auditing in Certificate Transparency. In *PETS*, 2022.
- [75] Marcela S. Melara, Aaron Blankstein, Joseph Bonneau, Edward W. Felten, and Michael J. Freedman. CONIKS: Bringing key transparency to end users. In *USENIX Security*, 2015.
- [76] Samir Jordan Menon and David J. Wu. Spiral: Fast, high-rate single-server PIR via FHE composition. In *S&P*, 2022.
- [77] Muhammad Haris Mughees, Hao Chen, and Ling Ren. OnionPIR: Response efficient single-server PIR. In *CCS*, 2021.
- [78] Moni Naor and Eylon Yogev. Bloom filters in adversarial environments. In *CRYPTO*, 2015.
- [79] Tatsuaki Okamoto and Shigenori Uchiyama. A new public-key cryptosystem as secure as factoring. In *EUROCRYPT*, 1998.
- [80] Rafail Ostrovsky and William E Skeith. A survey of single-database private information retrieval: Techniques and applications. In *PKC*, 2007.
- [81] Pascal Paillier. Public-key cryptosystems based on composite degree residuosity classes. In *EUROCRYPT*, 1999.
- [82] Bijeeta Pal, Mazharul Islam, Thomas Ristenpart, and Rahul Chatterjee. Might I Get Pwned: A second generation password breach alerting service. In *USENIX Security*, 2022.
- [83] Chris Peikert, Vinod Vaikuntanathan, and Brent Waters. A framework for efficient and composable oblivious transfer. In *CRYPTO*, 2008.
- [84] Giuseppe Persiano and Kevin Yeo. Limits of preprocessing for single-server PIR. In *SODA*, 2022.
- [85] Joel Reardon, Jeffrey Pound, and Ian Goldberg. Relational-complete private information retrieval. Technical report, University of Waterloo, CACR, 2007.
- [86] Oded Regev. On lattices, learning with errors, random linear codes, and cryptography. *Journal of the ACM*, 2009.
- [87] Sacha Servan-Schreiber, Kyle Hogan, and Srinivas Devadas. AdVeil: A private targeted-advertising ecosystem. Cryptology ePrint Archive, Report 2021/1032, 2021.
- [88] Elaine Shi, Waqar Aqeel, Balakrishnan Chandrasekaran, and Bruce Maggs. Puncturable pseudorandom sets and private information retrieval with near-optimal online bandwidth and time. In *CRYPTO*, 2021.
- [89] Emily Stark and Chris Thompson. Opt-in SCT auditing, 2020. <https://docs.google.com/document/d/1G1Jy8LJgSqJ-B673GnTYIG4b7XRw2ZLtvvS1rqFcl4A/edit>.
- [90] Volker Strassen. Gaussian elimination is not optimal. *Numerische mathematik*, 1969.
- [91] Kurt Thomas, Jennifer Pullman, Kevin Yeo, Ananth Raghunathan, Patrick Gage Kelley, Luca Invernizzi, Borbala Benko, Tadek Pietraszek, Sarvar Patel, Dan Boneh,

and Elie Bursztein. Protecting accounts from credential stuffing with password breach alerting. In *USENIX Security*, 2019.

- [92] Vish Viswanathan, Karthik Kumar, Thomas Willhalm, Blazej Filipiak Patrick Lu, and Sri Sakthivelu. Intel Memory Latency Checker. <https://www.intel.com/content/www/us/en/developer/articles/tool/intelr-memory-latency-checker.html>.
- [93] Frank Wang, Catherine Yun, Shafi Goldwasser, Vinod Vaikuntanathan, and Matei Zaharia. Splinter: Practical private queries on public data. In *NSDI*, 2017.
- [94] Ke Coby Wang and Michael K Reiter. Detecting stuffing of a user’s credentials at her own accounts. In *USENIX Security*, 2020.
- [95] Virginia Vassilevska Williams. Multiplying matrices faster than Coppersmith-Winograd. In *STOC*, 2012.
- [96] Zhenfei Zhang, Thomas Plantard, and Willy Susilo. On the CCA-1 security of somewhat homomorphic encryption over the integers. In Mark D. Ryan, Ben Smyth, and Guilin Wang, editors, *Information Security Practice and Experience*, 2012.
- [97] Mingxun Zhou, Wei-Kai Lin, Yiannis Tselekounis, and Elaine Shi. Optimal single-server private information retrieval. Cryptology ePrint Archive, Paper 2022/609, 2022.

## A Additional details on related work

For each PIR scheme from related work, we take its “optimal” entry size to be that for which the highest throughput was reported in the corresponding paper (or, if omitted, in a related paper). For each of our new PIR schemes (SimplePIR and DoublePIR), we compute its “optimal” entry size by executing the scheme on entries of increasing size, and selecting the entry size that yields the highest throughput. In Table 10, we display these entry sizes, along with each PIR scheme’s measured throughput on a database roughly 1 GB in size, with entries of the optimal size.

## B Client-state-recovery attacks on some existing PIR schemes

We demonstrate that several recent PIR schemes, including SealPIR and its descendants [6,8,9,76,77], are insecure against a certain type of active attack that enables the server to recover a client’s long-term, secret state. After such an attack, the server can learn the contents of all of the client’s PIR queries—even the ones made *before* the attacker compromised the server. In

	Database size		Max. achievable
	$N$	$\times d$	throughput/core
<b>Prior two-server PIR</b>			
DPF PIR [59]	$2^{25}$	32 B	5,381 MB/s*
XOR PIR	$2^{33}$	1 bit	6,067 MB/s*
XOR PIR fast	$2^{33}$	1 bit	11,797 MB/s*
<b>Prior single-server PIR</b>			
SealPIR [9]	$2^{22}$	288 B	97 MB/s
MulPIR [8]	$10^5$	40 KB	69 MB/s†
FastPIR [6]	$2^{20}$	1024 B	215 MB/s
OnionPIR [77]	$2^{15}$	30 KB	104 MB/s
Spiral [76]	$2^{14}$	100 KB	353 MB/s
SpiralPack [76]	$2^{15}$	30 KB	303 MB/s
SpiralStream [76]	$2^{15}$	30 KB	518 MB/s
SpiralStreamPack [76]	$2^{15}$	30 KB	1,314 MB/s
FrodoPIR [34]	$2^{20}$	1 KB	1,256 MB/s
<b>This work (single-server PIR)</b>			
<b>SimplePIR</b>	$2^{20}$	1 KB	10,305 MB/s
<b>DoublePIR</b>	$2^{33}$	1 bit	7,622 MB/s

Table 10: Maximal throughput measured for each PIR scheme, on databases of size roughly 1 GB, consisting of  $N$  entries each of size  $d$ . The entry sizes,  $d$ , are those for which the highest throughput was reported in the corresponding paper. \*The throughput is normalized by the number of servers, i.e., divided by two for 2-server PIR schemes. †We estimate MulPIR’s throughput from the measurements given in the paper, as no implementation is publicly available at this date.

contrast, our PIR schemes are not vulnerable to this type of attack (see Remark 3.2) and thus provide a form of forward secrecy (Definition B.1).

Fully reasoning about the security of single-server PIR schemes under active attack is nuanced and beyond the scope of this paper. We point out this active attack to highlight the fact that PIR schemes with secret long-term client state can carry additional security risks that our schemes do not.

The active attack applies to schemes in which the client holds and uses a persistent state across many queries. The state is a secret key for a homomorphic-encryption scheme. In particular, such schemes have the following syntax: first, the client runs a setup algorithm, which we call  $\text{Setup}_c$ . In contrast, in our new PIR schemes, it is the server who runs the setup algorithm. The  $\text{Setup}_c$  algorithm generates two pieces of information:

1. a public hint,  $\sigma_s$ , that the client transmits to the server, and
2. a private hint,  $\sigma_c$ , that the client keeps to itself.

Both  $\sigma_s$  and  $\sigma_c$  are re-used (by the server and the client respectively) across many queries: each time the client wants to query for some index  $j \in [N]$ , the client runs  $\text{Query}(\sigma_c, j) \rightarrow (\text{st}, \text{qu})$  and sends  $\text{qu}$  to the server. The server runs  $\text{Answer}(\text{db}, \sigma_s, \text{qu}) \rightarrow \text{ans}$  and sends  $\text{ans}$  to the client. Finally, the client runs  $\text{Recover}(\text{st}, \sigma_c, \text{ans}) \rightarrow d$  to recover the database record of interest. As we will show in the remainder of this section, this syntax can be problematic

as it lends itself to active attacks in which a malicious server recovers the client’s private, persistent state  $\sigma_c$ .

We note that the reason why SealPIR and related schemes reuse a persistent, secret client state over many queries is to save on communication. In these schemes,  $\sigma_c$  holds the client’s secret key for a homomorphic encryption scheme, while  $\sigma_s$  holds some “key-switching hints” associated with this key. The server uses these “key-switching hints” to perform query compression [9] or response packing [6].

**Definition B.1** (Forward secrecy against active attacks). Given a database size  $N \in \mathbb{N}$ , a number of queries  $Q \in \mathbb{N}$ , and a PIR scheme (Setup $_c$ , Query, Answer, Recover), consider the following experiment between a challenger and an adversary,

1. The challenger computes  $(\sigma_s, \sigma_c) \leftarrow \text{Setup}_c(1^\lambda)$  and sends  $\sigma_s$  to the adversary.
2. The adversary sends  $i_0, i_1 \in [N]$  to the challenger.
3. The challenger runs  $\text{Query}(\sigma_c, i_b) \rightarrow (\_, \text{qu})$ , and sends the query  $\text{qu}$  to the adversary.
4. For  $q \in \{1, \dots, Q\}$ :
  - The adversary chooses index  $j_q \in [N]$  and sends it to the challenger.
  - The challenger runs  $\text{Query}(\sigma_c, j_q) \rightarrow (\text{st}_q, \text{qu}_q)$  and sends the query  $\text{qu}_q$  to the adversary.
  - The adversary computes some (potentially malformed)  $\text{ans}'_q$  and sends it to the challenger.
  - The challenger runs  $\text{Recover}(\text{st}_q, \sigma_c, \text{ans}'_q) \rightarrow d$  and sends it to the adversary. (Note that  $d$  could be the failure symbol,  $\perp$ , if the Recover algorithm fails on  $\text{ans}'_q$ .)
5. The adversary outputs a guess  $\tilde{b} \in \{0, 1\}$ .

The PIR scheme satisfies  $(T, \epsilon)$ -forward secrecy against active attacks for  $Q$  queries if, for all  $N \in \mathbb{N}$  and all adversaries  $\mathcal{A}$  running in time at most  $T$  and making up to  $Q$  queries in the above experiment, it holds that

$$|\Pr[W_0] - \Pr[W_1]| \leq \epsilon,$$

where  $W_b$  denotes the event that  $\mathcal{A}$  outputs “1” in Experiment  $b$ , for  $b \in \{0, 1\}$ .

**Active attacks on some stateful PIR schemes.** Existing PIR schemes, which let the client reuse a secret key for homomorphic encryption indefinitely, do not provide forward secrecy against active attacks because their underlying homomorphic encryption schemes are not CCA1-secure [25, 43, 44, 70, 96]. In the remainder of this section, we show that there exist attacks that, with polynomially many PIR queries, recover the entire client state. For instance, consider SealPIR with  $d = 1$  (i.e., no recursion) and recall that the scheme uses BFV encryption [19, 42], where the secret keys are ring elements  $s = s_0 + s_1x + s_2x^2 + \dots + s_{n-1}x^{n-1}$  in the quotient ring  $R_q := \mathbb{Z}_q[x]/(x^n + 1)$  and  $p$  denotes the plaintext modulus. We modify the existing CCA attacks on BFV encryption [25]

to obtain Algorithm 1, which is an active attack on SealPIR. Conceptually, this attack recovers the client’s secret key by performing a binary search for each  $s_i \in \mathbb{Z}_q$ , for  $i \in [n]$ .

**Algorithm 1** Active attack for recovering  $s_i \in \mathbb{Z}_q$  in SealPIR

---

**Setup:** Let  $\mathcal{O}(\cdot) \leftarrow \text{Recover}(\text{st}, \sigma_c, \cdot)$ . (To implement  $\mathcal{O}(\cdot)$  in step 4 of the experiment in Definition B.1, first send  $\mathcal{C}$  any index in  $[N]$ , and then send  $\mathcal{C}$  the desired input to  $\mathcal{O}(\cdot)$ .)  
Define  $\Delta \leftarrow \lfloor q/p \rfloor$ .  
**Algorithm:** Let  $g \leftarrow \mathcal{O}(-1, 0)$ , where  $g \in \mathbb{Z}_q[x]/(x^n + 1)$ . Let  $t \leftarrow g[i]$  (where  $g[i]$  denotes the coefficient of  $x^i$  in  $g$ ).  
**if**  $t \neq 0$  **then**  
    Define  $M \leftarrow \Delta$ .  
**else**  
    Define  $M \leftarrow (q \bmod p) + \Delta$ .  
**end if**  
Perform binary search to find the smallest  $j \in \{0, 1, \dots, M\}$  such that  $\mathcal{O}(-1, j \cdot x^i)[i] = (t + 1) \bmod p$ .  
Output  $((t \cdot \Delta + \lfloor \frac{\Delta}{2} \rfloor + 1) - j) \bmod q$ .

---

*Analysis of Algorithm 1.* In SealPIR, when  $d = 1$ , the client’s private hint  $\sigma_c$  is the BFV secret key, i.e., a secret ring element  $s \in R_q$ , that the client uses to generate its queries. The server answers each client query with a single BFV ciphertext  $(a, b) \in R_q^2$ . Finally, it holds that

$$\begin{aligned} \mathcal{O}(-1, j \cdot x^i) &= \text{Dec}_{BFV}(s, (-1, j \cdot x^i)) \\ &= \text{Round}_\Delta(j \cdot x^i + s)/\Delta, \end{aligned} \quad (1)$$

where  $\text{Round}_\Delta(g)$  rounds each coefficient of the polynomial  $g$  to the closest multiple of  $\Delta := \lfloor q/p \rfloor$ .

By Equation (1), we see that  $t$  (defined as in Algorithm 1) is equal to  $\mathcal{O}(-1, 0)[i]$ , or, equivalently,  $\text{Round}_\Delta(s_i)/\Delta \in \mathbb{Z}_p$ . Therefore,  $s_i$  must fall into the following ranges:

- If  $t = 0$ , then

$$s_i \in \left\{ p \cdot \Delta - \left\lfloor \frac{\Delta}{2} \right\rfloor, \dots, q - 1 \right\} \cup \left\{ 0, \dots, \left\lfloor \frac{\Delta}{2} \right\rfloor \right\}.$$

- If  $t \neq 0$ , then  $s_i \in \{t \cdot \Delta - \lfloor \frac{\Delta}{2} \rfloor, \dots, t \cdot \Delta + \lfloor \frac{\Delta}{2} \rfloor\}$ .

Let  $j$  be the smallest value in  $\{0, \dots, M\}$  such that  $\text{Round}_\Delta(j + s_i)/\Delta = \text{Round}_\Delta(s_i)/\Delta + 1 \in \mathbb{Z}_p$ . The algorithm performs a binary search to find  $j$ , requiring roughly  $\log q$  queries to  $\mathcal{O}(\cdot)$ . We observe that, after this binary search,  $j$  will be exactly  $T - s_i$ , where  $T = t \cdot \Delta + \lfloor \frac{\Delta}{2} \rfloor + 1$  is the “next rounding boundary” of  $\text{Round}_\Delta(\cdot)$ . Finally, the algorithm outputs  $T - j$ , which is equal to  $s_i$ . Thus, it successfully recovers  $s_i \in \mathbb{Z}_q$ , with roughly  $\log q$  queries to  $\mathcal{O}(\cdot)$ .

To recover the entire secret key  $s$ , an attacker can perform the above binary search for all coefficients in  $\{s_i\}_{i \in [n]}$  in parallel. In other words, instead of asking the oracle  $\mathcal{O}$  for  $(-1, j \cdot x^i)$ , the attacker may ask for  $(-1, \sum_{i=1}^n j_i \cdot x^i)$ , where the  $\{j_i\}_{i \in [n]}$  are the query points checked by each of the  $n$  parallel binary searches. So, the total number of queries needed to recover  $s \in R_q$  remains roughly  $\log q$ . After this attack, the adversary

knows the client's secret key,  $s$ , and can thus learn all of the client's past and future PIR queries.

## C Additional material on SimplePIR

### C.1 Parameter selection

We instantiate SimplePIR with the parameter choices that maximize its throughput, while meeting the desired correctness and security requirements. Given a database size  $N$ , a correctness failure probability  $\delta$ , an adversary's runtime  $T$ , and a security failure probability  $\epsilon$ , we proceed as follows:

1. We first fix the ciphertext modulus  $q$  to be one of  $\{2^{16}, 2^{32}, 2^{64}\}$ , as modern hardware natively supports operations over  $\mathbb{Z}_q$  with these moduli.
2. We use LWE hardness estimates [7] to pick the LWE secret dimension  $n$  along with the LWE error distribution  $\chi$ , such that a collection of  $\sqrt{N}$  such LWE samples is  $(T, \epsilon)$ -secure against known attacks.
3. We compute the largest plaintext modulus,  $p$ , such that the chosen LWE parameters support at least  $\sqrt{N}$  homomorphic additions, with correctness error probability  $\delta$ .

In this work, we take  $\chi$  to be a discrete Gaussian distribution. We give concrete values for our parameters in Section 4.2.

### C.2 Correctness and security of SimplePIR

We now give the formal SimplePIR theorem statement:

**Theorem C.1** (SimplePIR). *On database size  $N \in \mathbb{N}$ , correctness failure probability  $\delta$ , adversary runtime  $T$ , and security failure probability  $\epsilon$ , let*

- $\chi$  be the discrete Gaussian distribution with variance  $\sigma^2$ ,
- $(n, q, \chi)$  be LWE parameters achieving  $(T, \epsilon)$ -security for  $\sqrt{N}$  LWE samples, and
- $p \in \mathbb{N}$  be a plaintext modulus chosen to satisfy

$$\lfloor q/p \rfloor \geq \sqrt{2} \cdot \sigma \cdot p \cdot N^{1/4} \cdot \sqrt{\ln(2/\delta)}. \quad (2)$$

Then, for a random LWE matrix  $\mathbf{A} \stackrel{\mathcal{R}}{\leftarrow} \mathbb{Z}_q^{\sqrt{N} \times n}$ , SimplePIR is a  $(T - O(\sqrt{N}), 2\epsilon)$ -secure PIR scheme on database size  $N$ , over plaintext space  $\mathbb{Z}_p$ , with correctness error  $\delta$ .

*Proof.* We prove correctness and security separately.

**Correctness.** Consider a client that interacts with the server to query for the database value at  $(i_{\text{row}}, i_{\text{col}}) \in [\sqrt{N}]^2$ . Let  $\mathbf{u}_{i_{\text{row}}}$  denote the unit vector  $i_{\text{row}}$  in  $\mathbb{Z}_q^{\sqrt{N}}$  (i.e., the vector of all zeros, with a single '1' at index  $i_{\text{row}}$ ). The Recover routine in

SimplePIR computes:

$$\begin{aligned} \hat{d} &= \text{ans}[i_{\text{row}}] - \text{hint}_c[i_{\text{row}}, :] \cdot \mathbf{s} \\ &= \mathbf{u}_{i_{\text{row}}}^T \cdot (\text{ans} - \text{hint}_c \cdot \mathbf{s}) \\ &= \mathbf{u}_{i_{\text{row}}}^T \cdot (\text{db} \cdot (\mathbf{A}\mathbf{s} + \mathbf{e} + \Delta \cdot \mathbf{u}_{i_{\text{col}}}) - (\text{db} \cdot \mathbf{A}) \cdot \mathbf{s}) \\ &= \mathbf{u}_{i_{\text{row}}}^T \cdot (\text{db} \cdot \mathbf{e} + \Delta \cdot \text{db} \cdot \mathbf{u}_{i_{\text{col}}}) \\ &= \text{db}[i_{\text{row}}, :] \cdot \mathbf{e} + \Delta \cdot \text{db}[i_{\text{row}}, i_{\text{col}}], \end{aligned}$$

where  $\text{db}[i_{\text{row}}, :]$  denotes row  $i_{\text{row}}$  of the database matrix  $\text{db}$ , and  $\text{db}[i_{\text{row}}, i_{\text{col}}]$  denotes the element at  $(i_{\text{row}}, i_{\text{col}})$  in  $\text{db}$ . Recovery succeeds when  $\text{Round}_\Delta(\hat{d})/\Delta$  is equal to  $\text{db}[i_{\text{row}}, i_{\text{col}}]$ . This happens if and only if  $|\text{db}[i_{\text{row}}, :] \cdot \mathbf{e}| < \Delta/2$ .

We can guarantee successful decryption by ensuring that  $\lfloor p/2 \rfloor \cdot \sqrt{N} \cdot |\mathbf{e}|_\infty < \Delta/2$  (because we can store the database entries—which are elements in  $\mathbb{Z}_p$ —as values in the range  $\{-\lfloor p/2 \rfloor, -\lfloor p/2 \rfloor + 1, \dots, \lfloor p/2 \rfloor - 1\}$ , so that they have maximal norm  $\lfloor p/2 \rfloor$ ). However, we instead use a tighter bound on  $|\text{db}[i_{\text{row}}, :] \cdot \mathbf{e}|$  by relying on properties of the error distribution,  $\chi$ . As a result, we will obtain a scheme in which decryption succeeds with probability  $1 - \delta$  (rather than 1) with a larger plaintext modulus,  $p$ .

Indeed, as  $\chi$  is the discrete Gaussian distribution with variance  $\sigma^2 = \frac{s^2}{2\pi}$  for some  $s > 0$ , by [68, Lemma 2.2][12, Lemma 2.4], for any  $T > 0$ , it holds that,

$$\Pr[|\text{db}[i_{\text{row}}, :] \cdot \mathbf{e}| \geq T \cdot s \cdot \|\text{db}[i_{\text{row}}, :]\|] < 2 \exp(-\pi \cdot T^2),$$

where  $\|\cdot\|$  denotes the Euclidean norm.

Taking  $T = \Delta/(2s \cdot \|\text{db}[i_{\text{row}}, :]\|)$ , we see that

$$\Pr[|\text{db}[i_{\text{row}}, :] \cdot \mathbf{e}| \geq \Delta/2] < \delta,$$

as long as

$$2 \exp\left(-\pi \cdot \left(\frac{\Delta}{2s \cdot \|\text{db}[i_{\text{row}}, :]\|}\right)^2\right) \leq \delta.$$

Equivalently, recovery fails with probability at most  $\delta$ , if

$$\Delta \geq 2s \cdot \|\text{db}[i_{\text{row}}, :]\| \cdot \sqrt{\frac{\ln(2/\delta)}{\pi}}.$$

Again, as we can store the elements in  $\text{db}$  in the range  $[-\lfloor p/2 \rfloor, \lfloor p/2 \rfloor]$ , we have that

$$\|\text{db}[i_{\text{row}}, :]\| \leq \sqrt{\sqrt{N} \cdot \lfloor p/2 \rfloor^2} = N^{1/4} \cdot \lfloor p/2 \rfloor.$$

In addition, we know that  $\Delta = \lfloor q/p \rfloor$  and  $s = \sigma \cdot \sqrt{2\pi}$ . Thus, recovery fails with probability at most  $\delta$ , as long as:

$$\lfloor q/p \rfloor \geq \sigma \cdot \sqrt{2\pi} \cdot p \cdot N^{1/4} \cdot \sqrt{\frac{\ln(2/\delta)}{\pi}},$$

or, equivalently,

$$\lfloor q/p \rfloor \geq \sqrt{2} \cdot \sigma \cdot p \cdot N^{1/4} \cdot \sqrt{\ln(2/\delta)},$$

By Equation (2), this condition holds. Thus, SimplePIR is correct, except with error probability at most  $\delta$ .

**Security.** Security follows from the fact that the LWE problem is hard, even when the LWE matrix  $\mathbf{A}$  is reused across many independent trials, each using an independently generated LWE secret  $\mathbf{s} \leftarrow \mathbb{Z}_q^n$ , as shown in prior work [83, Lemma 7.3]. We connect the security of SimplePIR to the hardness of LWE with the following lemma:

**Lemma C.2.** *Let  $N \in \mathbb{N}$  be the database size,  $(n, q, \chi)$  be the LWE parameters, and  $\mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n}$  be the random LWE matrix used in SimplePIR. For any  $i \in [N]$ , we define the distribution*

$$\mathcal{Q}_i = \{(\mathbf{A}, \text{qu}_i) : \_, \text{qu}_i \leftarrow \text{Query}(i)\}.$$

*If the  $(n, q, \chi)$ -LWE problem with  $\sqrt{N}$  samples is  $(T, \epsilon)$ -hard, then any algorithm running in time  $T - O(\sqrt{N})$  can have success probability at most  $\epsilon$  in distinguishing  $\mathcal{Q}_i$  from the distribution  $\{(\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow \mathbb{Z}_q^{\sqrt{N}}\}$ .*

*Proof.* Consider any index  $i \in [N]$ . We decompose  $i$  into the pair of coordinates  $(i_{\text{row}}, i_{\text{col}}) \in [\sqrt{N}]^2$ , as done in the Query routine in SimplePIR, and let  $\mathbf{u}_{i_{\text{col}}}$  denote unit vector  $i_{\text{col}}$  in  $\mathbb{Z}_q^{\sqrt{N}}$ . Additionally, we define the following distributions:

- $\mathcal{D}_1 = \{(\mathbf{A}, \mathbf{A}\mathbf{s} + \mathbf{e}) : \mathbf{s} \leftarrow \mathbb{Z}_q^n, \mathbf{e} \leftarrow \chi^{\sqrt{N}}\}$ , and
- $\mathcal{D}_2 = \{(\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow \mathbb{Z}_q^{\sqrt{N}}\}$ .

Now, consider the simulator  $\mathcal{S}$  that, given as input  $(\mathbf{A}, \mathbf{v}) \in \mathbb{Z}_q^{\sqrt{N} \times n} \times \mathbb{Z}_q^{\sqrt{N}}$ , computes and outputs  $(\mathbf{A}, \mathbf{v} + \lfloor q/p \rfloor \cdot \mathbf{u}_{i_{\text{col}}})$ . We observe that:

- when  $(\mathbf{A}, \mathbf{v})$  is sampled from  $\mathcal{D}_1$ , the simulator's output is distributed identically to  $\mathcal{Q}_i$ .
- when  $(\mathbf{A}, \mathbf{v})$  is sampled from  $\mathcal{D}_2$ , the simulator's output is distributed identically to  $\mathcal{D}_2$ .

As the  $(n, q, \chi)$ -LWE problem with  $\sqrt{N}$  samples is  $(T, \epsilon)$ -hard, we know that any algorithm running in time  $T$  has advantage at most  $\epsilon$  in distinguishing between  $\mathcal{D}_1$  and  $\mathcal{D}_2$ . Our simulator  $\mathcal{S}$  runs in time  $O(\sqrt{N})$ . Thus, it must hold that any algorithm distinguishing between  $\mathcal{Q}_i$  and  $\mathcal{D}_2$  in time at most  $T - O(\sqrt{N})$  can have success probability at most  $\epsilon$ .  $\square$

Lemma C.2 shows that any algorithm running in time  $T - O(\sqrt{N})$  can distinguish the queries made by a client in SimplePIR from random vectors in  $\mathbb{Z}_q^{\sqrt{N}}$  with success probability at most  $\epsilon$ . Therefore, by the triangle inequality, any algorithm running in time  $T - O(\sqrt{N})$  can distinguish the queries made by a client in SimplePIR to any pair of indices  $i \in [N]$  and  $j \in [N]$  with success probability at most  $2\epsilon$ . In other words, SimplePIR is  $(T - O(\sqrt{N}), 2\epsilon)$ -secure.  $\square$

We can extend the security definition in Section 3.2 to handle  $Q$  queries by requiring that any two sequences of  $Q$  queries produce indistinguishable query distributions. A hybrid argument proves the following corollary:

**Corollary C.3** ( *$Q$ -query security of SimplePIR*). *For any sequence  $I \in [N]^Q$  of  $Q$  indices, we define  $\mathcal{D}_I$  to be the query distribution that it induces, i.e.,*

$$\mathcal{D}_I = \{(\mathbf{A}, \text{qu}_k) : \_, \text{qu}_k \leftarrow \text{Query}(I_k)\}_{k \in [Q]}$$

*Also, we define the random query distribution,  $\mathcal{R}$ :*

$$\mathcal{R} = \left\{ (\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow \mathbb{Z}_q^{\sqrt{N}} \right\}_{k \in [Q]}$$

*If the  $(n, q, \chi)$ -LWE problem with  $\sqrt{N}$  samples is  $(T, \epsilon)$ -hard, then  $\mathcal{D}_I$  is  $(T - O(Qn\sqrt{N}), Q\epsilon)$ -indistinguishable from  $\mathcal{R}$ .*

*Proof.* The proof follows by a standard hybrid argument. Let  $I \in [N]^Q$  be a sequence of  $Q$  indices. For any  $j \in [Q]$ , we define the hybrid distribution,  $H_j$ :

$$H_j = \left\{ (\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow \mathbb{Z}_q^{\sqrt{N}} \right\}_{k \in \{1, \dots, j\}} \\ \cup \{(\mathbf{A}, \text{qu}_k) : \_, \text{qu}_k \leftarrow \text{Query}(I_k)\}_{k \in \{j+1, \dots, Q\}}.$$

In other words,  $H_j$  is the distribution where the first  $j$  queries are truly random vectors, and the remaining  $(Q - j)$  queries are sampled from  $\mathcal{D}_I$ . As such,  $H_0 = \mathcal{D}_I$ , while  $H_Q = \mathcal{R}$ .

Now, consider any algorithm  $\mathcal{A}$  that runs in time  $t$  and distinguishes between distributions  $H_0$  and  $H_Q$  with advantage at least  $p$ . For any  $j \in [Q]$ , we define  $p_j$  to be the probability that  $\mathcal{A}$  outputs 1 when given samples from  $H_j$ . By definition,

$$p = |p_Q - p_0| = \left| \sum_{j=1}^Q (p_j - p_{j-1}) \right|.$$

Thus, there exists some  $j^* \in [Q]$  such that

$$|p_{j^*} - p_{j^*-1}| \geq p/Q.$$

Using  $\mathcal{A}$  as a subroutine, we then build an algorithm  $\mathcal{B}$  that distinguishes between the distributions  $\mathcal{Q}_{I_{j^*}}$  (defined in Lemma C.2) and  $\{(\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow \mathbb{Z}_q^{\sqrt{N}}\}$ .

---

**Algorithm 2**  $\mathcal{B}$ , on input  $(\mathbf{A}, \mathbf{b})$ :

---

```

for  $j = 1, \dots, j^* - 1$  do
  Compute  $\mathbf{r}_j \leftarrow \mathbb{Z}_q^{\sqrt{N}}$ 
end for
for  $j = j^* + 1, \dots, Q$  do
  Compute  $\text{qu}_j \leftarrow \text{Query}(I_j)$ 
end for
Output

```

$$\mathcal{A} \left( (\mathbf{A}, \mathbf{r}_1), \dots, (\mathbf{A}, \mathbf{r}_{j^*-1}), (\mathbf{A}, \mathbf{b}), (\mathbf{A}, \text{qu}_{j^*+1}), \dots, (\mathbf{A}, \text{qu}_Q) \right)$$


---

Then, we observe that:

- When  $\mathcal{B}$  is given an input sampled from  $\mathcal{Q}_{I_{j^*}}$ , the input that  $\mathcal{B}$  feeds to  $\mathcal{A}$  is distributed following  $H_{j^*}$ . As such,  $\mathcal{B}$  outputs 1 with probability  $p_{j^*}$ .

- When  $\mathcal{B}$  is given an input sampled from  $\{(\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow^{\mathcal{R}} \mathbb{Z}_q^{\sqrt{N}}\}$ , the input that  $\mathcal{B}$  feeds to  $\mathcal{A}$  is distributed following  $H_{j^*+1}$ . As such,  $\mathcal{B}$  outputs 1 with probability  $p_{j^*+1}$ .

So,  $\mathcal{B}$  distinguishes between the distributions  $\mathcal{Q}_{I_{j^*}}$  and  $\{(\mathbf{A}, \mathbf{r}) : \mathbf{r} \leftarrow^{\mathcal{R}} \mathbb{Z}_q^{\sqrt{N}}\}$  with advantage at least  $|p_{j^*} - p_{j^*-1}| \geq p/Q$ . Further,  $\mathcal{B}$  runs in time  $t + O(Qn\sqrt{N})$ , as Query takes time  $O(n\sqrt{N})$ .

By Lemma C.2, we conclude that, if the  $(n, q, \chi)$ -LWE problem with  $\sqrt{N}$  samples is  $(T, \epsilon)$ -hard and  $t + O(Qn\sqrt{N}) \leq T - O(\sqrt{N})$ , it must hold that  $p/Q < \epsilon$ . In other words, if the  $(n, q, \chi)$ -LWE problem with  $\sqrt{N}$  samples is  $(T, \epsilon)$ -hard, then  $\mathcal{D}_I$  is  $(O(T - Qn\sqrt{N}), Q\epsilon)$ -indistinguishable from  $\mathcal{R}$ .  $\square$

### C.3 Additional extensions and optimizations

**Answering many client queries at once.** To answer a client query, the SimplePIR server multiplies the  $\sqrt{N}$ -by- $\sqrt{N}$  database matrix by the client’s dimension- $\sqrt{N}$  query vector, in  $2N$  operations. If the server wants to answer  $\sqrt{N}$  client queries at once, it can use fast matrix-multiplication algorithms [29, 95] to compute this more efficiently than answering each query independently—in  $o(N^{3/2})$  time. This observation comes directly from prior work [13, 71]. We have not yet experimented with this optimization, since it is not clear to us that the asymptotic efficiency gain will translate to concrete cost improvements.

**Faster preprocessing.** In its preprocessing step, the server computes the matrix-matrix product of the database with the LWE matrix,  $\mathbf{A}$ . The asymptotic performance of this step can be improved using fast matrix-multiplication algorithms [29, 65, 90, 95]. We suspect that such a refinement will not improve the concrete performance for our parameter sizes.

**Handling database updates.** As discussed in Remark 3.1, if the database changes, the server inherently needs perform some of its preprocessing work again. In SimplePIR, the amount of preprocessing work that needs to be repeated is proportional to the number of rows in the database matrix whose contents have changed. More specifically, the preprocessing phase in SimplePIR computes  $\text{hint}_c$  to be the product of the database matrix,  $\text{db}$ , and the LWE matrix,  $\mathbf{A}$ . Thus, when some row  $\mathbf{d}$  of the database changes, only the corresponding row in  $\text{hint}_c$  changes and must be updated. In other words, it is sufficient for the server to compute  $\mathbf{d} \cdot \mathbf{A} \in \mathbb{Z}_q^{1 \times n}$  (rather than the much more expensive  $\text{db} \cdot \mathbf{A} \in \mathbb{Z}_q^{\sqrt{N} \times n}$ ) and send it back to the client. Additions and deletions of rows in  $\text{db}$  can be handled similarly.

**Decreasing the online download with local rounding.** In SimplePIR’s online phase, the client downloads a vector of  $\sqrt{N}$  elements in the ciphertext space,  $\mathbb{Z}_q$ , performs a linear computation on them, and finally rounds the result to map it into the message space,  $\mathbb{Z}_p$ . Prior work [38, Lemma 3.2][39, Lemma 4.1][18, Lemma 1] observes that, for any value  $v \in \mathbb{Z}_q$  that is

“close” to a multiple of  $\Delta$ , and for a random  $r \in \mathbb{Z}_q$ , it holds with high probability that

$$\text{Round}_{\Delta}(v) = \text{Round}_{\Delta}(v + r) - \text{Round}_{\Delta}(r).$$

Using this observation (with  $v \leftarrow \text{ans} - \text{hint}_c \cdot s$  and  $r \leftarrow \text{hint}_c \cdot s$ ), the server and the client could each locally round the values they contribute to the client’s final computation ( $\text{ans}$  and  $\text{hint}_c \cdot s$ , respectively), decreasing the online download to only  $\sqrt{N}$  elements in  $\mathbb{Z}_p$ . If the LWE parameters are chosen appropriately (namely,  $\lfloor q/p \rfloor \geq \sqrt{2} \cdot \frac{4}{\delta} \cdot \sigma \cdot p \cdot N^{1/4} \cdot \sqrt{\ln(8/\delta)}$ ), then this local rounding does not affect the scheme’s correctness (i.e., recovery still succeeds, except with probability  $\delta$ ).

In practice, this optimization decreases the online download by a factor of  $\frac{\log q}{\log p} \approx 3\times$ , at the expense of requiring a much smaller plaintext modulus,  $p$ , to maintain correctness and thus markedly reducing SimplePIR’s throughput. Exploring other approaches to local rounding [20] or to correcting the errors that it introduces are interesting directions for future work.

## D Linearly homomorphic encryption with preprocessing

In this section, we abstract out the key property of Regev encryption that enables SimplePIR’s high throughput. We expect this primitive to have applications beyond PIR.

Specifically, we introduce the notion of *linearly homomorphic encryption with preprocessing*. This notion is similar to standard linearly homomorphic encryption: given an encryption of vector  $\mathbf{v} \in \mathbb{Z}_p^m$  (where the ring  $\mathbb{Z}_p$  and dimension  $m$  are parameters of the scheme), it is possible to homomorphically evaluate any linear function  $f : \mathbb{Z}_p^m \rightarrow \mathbb{Z}_p$  on the encrypted values; after evaluation, the resulting ciphertext decrypts to the function value  $f(v)$ . What distinguishes our use of Regev encryption from standard linearly homomorphic encryption is that we can preprocess the function  $f$  to speed up the homomorphic-evaluation operation. In particular, given a linear function  $f$ , we compute a preprocessed hint,  $\text{hint}_f$ . The hint enables homomorphic evaluation of  $f$  on ciphertexts encrypting any  $\mathbf{v} \in \mathbb{Z}_p^m$  at very low computational cost—nearly the cost of evaluating  $f$  on a plaintext vector.

### D.1 Definition

Formally, a *linearly homomorphic encryption scheme with preprocessing* consists of the following efficient algorithms, which are implicitly parameterized by a security parameter  $n \in \mathbb{N}$ , a correctness parameter  $\delta \in \mathbb{N}$ , a plaintext dimension  $m \in \mathbb{N}$ , a key space  $\mathcal{K}$ , and a plaintext space  $\mathbb{Z}_p$ :

- Setup()  $\rightarrow$  pp. Sample and output the public parameters.
- Enc(pp, sk,  $\mathbf{v}$ )  $\rightarrow$  ct. Given the public parameters pp, a secret key  $\text{sk} \in \mathcal{K}$ , and a vector  $\mathbf{v} \in \mathbb{Z}_p^m$ , output a ciphertext ct.
- Preproc(pp,  $f$ )  $\rightarrow$   $\text{hint}_f$ . Given public parameters pp and a linear function  $f : \mathbb{Z}_p^m \rightarrow \mathbb{Z}_p$ , output a function hint  $\text{hint}_f$ .

Apply(pp,  $f$ , ct)  $\rightarrow$  ct $_f$ . Given the public parameters pp, a linear function  $f: \mathbb{Z}_p^m \rightarrow \mathbb{Z}_p$ , and a ciphertext ct, produce a ciphertext ct $_f$  that encrypts the value of the function  $f$  applied to the vector encrypted by ct.

Dec(sk, hint $_f$ , ct $_f$ )  $\rightarrow$   $d$ . Given a secret key sk, the function hint hint $_f$ , and a ciphertext ct $_f$ , output a decrypted value  $d$ .

**Correctness.** After encryption, preprocessing and application of the linear function, and decryption are carried out sequentially, the output should be the correct function value, except with negligible probability in the implicit correctness parameter. Formally, we say that the encryption scheme has *correctness error*  $\delta$  if, for all  $sk \in \mathcal{K}$ , for all  $\mathbf{v} \in \mathbb{Z}_p^m$ , and linear functions  $f \in \mathbb{Z}_p^m \rightarrow \mathbb{Z}_p$ , the following probability is at least  $1 - \delta$ :

$$\Pr \left[ \begin{array}{l} pp \leftarrow \text{Setup}() \\ ct \leftarrow \text{Enc}(pp, sk, \mathbf{v}) \\ \text{hint}_f \leftarrow \text{Preproc}(pp, f) \\ ct_f \leftarrow \text{Apply}(pp, f, ct) \\ d \leftarrow \text{Dec}(sk, \text{hint}_f, ct_f) \end{array} \middle| d = f(\mathbf{v}) \right].$$

**Security.** We require semantic security [51]. In other words, ct should reveal no information about the encrypted vector  $\mathbf{v} = (v_1, \dots, v_m)$ . Formally, given vectors  $\mathbf{v}_0, \mathbf{v}_1 \in \mathbb{Z}_p^m$ , for  $b \in \{0, 1\}$ , we define the value  $W_b$  as:

$$W_b := \Pr \left[ \begin{array}{l} \mathcal{A}(pp, ct) = 1: \\ pp \leftarrow \text{Setup}() \\ sk \leftarrow \mathcal{K} \\ ct \leftarrow \text{Enc}(pp, sk, \mathbf{v}_b) \end{array} \right].$$

Then, we say that the encryption scheme is *secure* if, for all efficient adversaries  $\mathcal{A}$ , for all  $\mathbf{v}_0, \mathbf{v}_1 \in \mathbb{Z}_p^m$ , it holds that  $|W_0 - W_1|$  is negligible in the implicit security parameter.

## D.2 Construction

In Figure 12, we construct a linear homomorphic encryption scheme with preprocessing from Regev encryption. The construction yields the following theorem:

**Theorem D.1.** *Assume the  $(n, q, \chi)$ -LWE problem with  $m$  samples is  $(T, \epsilon)$ -hard, and let  $p \in \mathbb{N}$  be a suitable plaintext modulus<sup>3</sup> for Regev encryption with these parameters supporting  $m$  homomorphic additions. Then, the construction of Figure 12 is a  $(T - O(m), 2\epsilon)$ -secure linearly homomorphic encryption scheme with preprocessing for plaintext dimension  $m$  and plaintext modulus  $p$  where:*

- the hint consists of  $n$  elements of  $\mathbb{Z}_q$ ,
- the ciphertext encrypting a vector in  $\mathbb{Z}_p^m$  consists of  $m$  elements of  $\mathbb{Z}_q$ ,
- the evaluation routine Apply requires  $m$  additions and  $m$  multiplications in  $\mathbb{Z}_q$ , and
- the ciphertext that Apply outputs consists of one  $\mathbb{Z}_q$  element.

<sup>3</sup>To support arbitrary linear functions, it is necessary to take  $p|q$ . To support linear functions whose output will not wrap around  $\bmod p$  (as is the case for PIR), such a restriction on  $p$  is not necessary.

	Size			Running time (per bit)			
	Hint	Key	Ct/bit	Preproc	Enc	Apply	Dec
Factoring [79, 81]	n/a	$\lambda^5$	$\lambda^2$	n/a	$\lambda^2$	$\lambda^2$	$\lambda^2$
EC ElGamal [40]	n/a	$\lambda$	$\lambda$	n/a	$\lambda^2$	$\lambda$	$\lambda$
Regev [83, 86]	n/a	$\lambda$	1	n/a	$\lambda$	$\lambda$	$\lambda$
Ring LWE [72]	n/a	$\lambda$	1	n/a	1	1	1
<b>This work (LWE)</b>	$\lambda$	$\lambda$	1	$\lambda$	$\lambda$	1	$\lambda$

Table 11: Comparison of linearly homomorphic encryption schemes on plaintext dimension  $m$ . We suppress  $\log \lambda$  factors and we assume that it is possible to perform modular multiplication in linear time. Since our scheme uses plain LWE instead of ring LWE, it is simpler to implement—there is no need for polynomial arithmetic, fast Fourier transforms, or a modulus of special form.

It is worth emphasizing the efficiency of the construction implied by Theorem D.1. By our choice of LWE parameters in Equation (2), the number of bits encoding each ciphertext element (i.e.,  $\log q$ ) grows logarithmically with the plaintext dimension  $m$  and is independent of the security parameter  $n$ . Thus, our construction has:

- a hint size that grows only logarithmically with the size of the linear function  $f$ ,
- ciphertext size that is only a  $\log m$  factor larger than the corresponding plaintext (and that is independent of the security parameter), and
- a homomorphic evaluation routine that is essentially as fast as plaintext evaluation of  $f$ , as homomorphic evaluation just requires computing  $f$  over  $\mathbb{Z}_q$  instead of over  $\mathbb{Z}_p$ . This evaluation time is independent of the security parameter.

In contrast, when using Paillier [81] or ElGamal-based [40] linearly homomorphic encryption, the ciphertext size and evaluation time grow quadratically or cubically with the security parameter (Table 11).

## D.3 Abstract view of SimplePIR

In Figure 13, we re-express the SimplePIR construction of Section 4 in the language of linear homomorphic encryption with preprocessing. The PIR scheme of Figure 13 demonstrates that, given a linearly homomorphic encryption scheme with preprocessing, there exists a PIR scheme with  $\sqrt{N}$  communication cost and for which the server’s online computation time consists of running the encryption’s Apply algorithm  $\sqrt{N}$  times, each on a vector of dimension  $\sqrt{N}$ . The client’s hint—which it must fetch in an offline phase—consists of  $\sqrt{N}$  hints for the encryption scheme, each computed for a linear function with  $\sqrt{N}$  inputs.

An interesting task for future work would be to construct a linear homomorphic encryption scheme with preprocessing with smaller hints. This would reduce the offline communication in the resulting PIR schemes.

*Parameters:* LWE parameters  $(n, q, \chi)$ , a plaintext modulus  $p \ll q$ . Define  $\Delta := \lfloor q/p \rfloor$ . The keyspace is  $\mathcal{K} = \mathbb{Z}_q^n$ .

Setup()  $\rightarrow$  pp  $\in \mathbb{Z}_q^{m \times n}$

- Sample  $\mathbf{A} \leftarrow^{\mathcal{R}} \mathbb{Z}_q^{m \times n}$ .
- Return pp  $\leftarrow \mathbf{A}$ .

Enc(pp, sk,  $\mathbf{v}$ )  $\rightarrow$  ct  $\in \mathbb{Z}_q^m$

- Parse  $(\mathbf{A} \in \mathbb{Z}_q^{m \times n}) \leftarrow$  pp.
- Sample  $\mathbf{e} \leftarrow^{\mathcal{R}} \chi^m$ .
- Return ct  $\leftarrow \mathbf{A} \cdot \text{sk} + \mathbf{e} + \Delta \cdot \mathbf{v} \in \mathbb{Z}_q^m$ .

Preproc(pp,  $f$ )  $\rightarrow$  hint $_f \in \mathbb{Z}_q^{1 \times n}$

- Parse  $(\mathbf{A} \in \mathbb{Z}_q^{m \times n}) \leftarrow$  pp.
- Parse  $(\mathbf{f} \in \mathbb{Z}_q^{1 \times m}) \leftarrow f$  as a row vector.
- Return hint $_f \leftarrow \mathbf{f} \cdot \mathbf{A}$ .

Apply(pp,  $f$ , ct)  $\rightarrow$  ct $_f \in \mathbb{Z}_q$

- Parse:
  - $(\mathbf{f} \in \mathbb{Z}_q^{1 \times m}) \leftarrow f$  as a row vector,
  - $(\mathbf{c} \in \mathbb{Z}_q^m) \leftarrow$  ct.
- Return ct $_f \leftarrow \mathbf{f} \cdot \mathbf{c} \in \mathbb{Z}_q$ .

Dec(pp, sk, hint $_f$ , ct $_f$ )  $\rightarrow$   $d \in \mathbb{Z}_p$

- Compute  $\hat{d} \leftarrow \text{ct}_f - \text{hint}_f \cdot \text{sk} \in \mathbb{Z}_q$ .
- Let  $v \in \mathbb{Z}_q$  be Round $_{\Delta}(\hat{d})$ , which is  $\hat{d}$  rounded to the nearest integral multiple of  $\Delta$ .
- Return  $d \leftarrow v/\Delta \in \mathbb{Z}_p$ .

Figure 12: Reggev Encryption expressed as linear homomorphic encryption with preprocessing.

## E Additional material on DoublePIR

We give a formal description of DoublePIR in Figure 14.

### E.1 Parameter selection

As for SimplePIR, we take the ciphertext modulus,  $q$ , to be  $2^{32}$  and the LWE distribution,  $\chi$ , to be a discrete Gaussian in DoublePIR. We use hardness estimates [7] to select appropriate choices of  $\chi$  and of the LWE dimension,  $n$ . Our selection of the plaintext modulus,  $p$ , in DoublePIR is slightly more conservative than that for SimplePIR, as for correctness to hold, the recovery routine must succeed for:

- the single element read from the first-level database, and
- the  $(n + 1)$  elements read from the second-level database.

We detail our calculation of  $p$  in Theorem E.1. Here, we give sample choices of the plaintext modulus,  $p$ , for different database sizes:

Database size $N$ :	$2^{26}$	$2^{28}$	$2^{30}$	$2^{34}$	$2^{38}$	$2^{42}$
Plaintext modulus $p$ :	929	781	657	464	328	231

*Parameters:* a database size  $N$ , a linearly homomorphic encryption with preprocessing scheme (Setup, Enc, Preproc, Apply, Dec) with plaintext dimension  $\sqrt{N}$ , plaintext space  $\mathbb{Z}_p$ , key space  $\mathcal{K}$ , and public parameters pp, computed as pp  $\leftarrow$  Setup(). Our construction uses the following helper routine:

ToFuncs(db)  $\rightarrow (f_1, \dots, f_{\sqrt{N}})$ :

- View the database db as vectors  $\mathbf{d}_1, \dots, \mathbf{d}_{\sqrt{N}} \in \mathbb{Z}_p^{\sqrt{N}}$ .
- For all  $i \in [\sqrt{N}]$ , define  $f_i(\mathbf{x}) := \langle \mathbf{d}_i, \mathbf{x} \rangle \in \mathbb{Z}_p$ . (Here,  $f_i$  is a linear function on  $\mathbb{Z}_p^{\sqrt{N}}$  to  $\mathbb{Z}_p$ .)
- Output  $(f_1, \dots, f_{\sqrt{N}})$ .

Setup(db)  $\rightarrow$  (hint $_s$ , hint $_c$ ).

- Let  $(f_1, \dots, f_{\sqrt{N}}) \leftarrow$  ToFuncs(db).
- For all  $i \in [\sqrt{N}]$ , compute: hint $_i \leftarrow$  Preproc(pp,  $f_i$ ).
- Set hint $_c \leftarrow$  (hint $_1, \dots, \text{hint}_{\sqrt{N}}$ ).
- Return  $(\_, \text{hint}_c)$ .

Query( $i$ )  $\rightarrow$  (st, qu).

- Write  $i$  as a pair  $(j, k) \in [\sqrt{N}] \times [\sqrt{N}]$ .
- Sample sk  $\leftarrow^{\mathcal{R}} \mathcal{K}$ .
- Compute qu  $\leftarrow$  Enc(pp, sk,  $\mathbf{u}_j \in \mathbb{Z}_p^{\sqrt{N}}$ ), where  $\mathbf{u}_j$  is the vector of all zeros with a single ‘1’ at index  $j$ .
- Set st  $\leftarrow$  (sk,  $k$ ) and return (st, qu).

Answer(db, hint $_s$ , qu)  $\rightarrow$  ans.

- Let  $(f_1, \dots, f_{\sqrt{N}}) \leftarrow$  ToFuncs(db).
- For all  $i \in [\sqrt{N}]$ , compute: ans $_i \leftarrow$  Apply(pp,  $f_i$ , qu).
- Return ans  $\leftarrow$  (ans $_1, \dots, \text{ans}_{\sqrt{N}}$ ).

Recover(st, hint $_c$ , ans)  $\rightarrow$   $d \in \mathbb{Z}_p$ .

- Parse (sk,  $k$ )  $\leftarrow$  st.
- Parse (hint $_1, \dots, \text{hint}_{\sqrt{N}}) \leftarrow$  hint $_c$ .
- Parse (ct $_1, \dots, \text{ct}_{\sqrt{N}}) \leftarrow$  ans.
- Return  $d \leftarrow$  Dec(pp, sk, hint $_k$ , ct $_k$ ).

Figure 13: The SimplePIR protocol, expressed in terms of linearly homomorphic encryption with preprocessing.



**Construction: DoublePIR.** The parameters of the construction are a database size  $N$ , database dimensions  $\ell$  and  $m$  (such that  $\ell \cdot m \geq N$ ), LWE parameters  $(n, q, \chi)$ , plaintext modulus  $p \ll q$ , and LWE matrices  $\mathbf{A}_1 \in \mathbb{Z}_q^{m \times n}$  and  $\mathbf{A}_2 \in \mathbb{Z}_q^{\ell \times n}$  (sampled in practice using a hash function). The database consists of  $N$  values in  $\mathbb{Z}_p$ , represented as a matrix in  $\mathbb{Z}_p^{\ell \times m}$ . We define the scalars  $\kappa := \left\lceil \frac{\log q}{\log p} \right\rceil$  and  $\Delta := \left\lfloor \frac{q}{p} \right\rfloor$ .

We use the helper routines:

- **Decomp:**  $\mathbb{Z}_q^{a \times b} \rightarrow \mathbb{Z}_q^{\kappa a \times b}$ , which writes each  $\mathbb{Z}_q$  element as its base- $p$  decomposition,
- **Recomp:**  $\mathbb{Z}_q^{\kappa a \times b} \rightarrow \mathbb{Z}_q^{a \times b}$ , which interprets each  $\kappa \times 1$  submatrix of its input as a base- $p$  decomposition of a  $\mathbb{Z}_q$  element and outputs the matrix of these elements.
- **Round $_{\Delta}$ :**  $\mathbb{Z}_q^a \rightarrow \mathbb{Z}_q^a$ , which rounds each entry to the nearest integral multiple of  $\Delta$ .

Setup( $\text{db} \in \mathbb{Z}_p^{\ell \times m}$ )  $\rightarrow$  ( $\text{hint}_s, \text{hint}_c$ ).

- Compute  $\begin{cases} \text{hint}_s \leftarrow \text{Decomp}(\mathbf{A}_1^T \cdot \text{db}^T) \in \mathbb{Z}_q^{\kappa n \times \ell} \\ \text{hint}_c \leftarrow \text{hint}_s \cdot \mathbf{A}_2 \in \mathbb{Z}_q^{\kappa n \times n} \end{cases}$
- Return ( $\text{hint}_s, \text{hint}_c$ ).

Query( $i \in [N]$ )  $\rightarrow$  ( $\text{st}, \text{qu}$ ).

- Write  $i$  as a pair  $(i_{\text{row}}, i_{\text{col}}) \in [\ell] \times [m]$ .
- Sample  $\begin{cases} \mathbf{s}_1 \stackrel{\mathbb{R}}{\leftarrow} \mathbb{Z}_q^n, & \mathbf{e}_1 \stackrel{\mathbb{R}}{\leftarrow} \chi^m \in \mathbb{Z}_q^m, \\ \mathbf{s}_2 \stackrel{\mathbb{R}}{\leftarrow} \mathbb{Z}_q^n, & \mathbf{e}_2 \stackrel{\mathbb{R}}{\leftarrow} \chi^\ell \in \mathbb{Z}_q^\ell. \end{cases}$
- Compute  $\begin{cases} \mathbf{c}_1 \leftarrow (\mathbf{A}_1 \cdot \mathbf{s}_1 + \mathbf{e}_1 + \Delta \cdot \mathbf{u}_{i_{\text{col}}}) \in \mathbb{Z}_q^m, \\ \mathbf{c}_2 \leftarrow (\mathbf{A}_2 \cdot \mathbf{s}_2 + \mathbf{e}_2 + \Delta \cdot \mathbf{u}_{i_{\text{row}}}) \in \mathbb{Z}_q^\ell, \end{cases}$   
where  $\mathbf{u}_{i^*}$  is the vector of all zeros with a single ‘1’ at index  $i^*$ .
- Return ( $\text{st}, \text{qu}$ )  $\leftarrow ((\mathbf{s}_1, \mathbf{s}_2), (\mathbf{c}_1, \mathbf{c}_2))$ .

Answer( $\text{db}, \text{hint}_s, \text{qu}$ )  $\rightarrow$  ans.

- Parse  $(\mathbf{c}_1 \in \mathbb{Z}_q^m, \mathbf{c}_2 \in \mathbb{Z}_q^\ell) \leftarrow \text{qu}$ .
- Compute  $\begin{cases} \text{ans}_1 \leftarrow \text{Decomp}(\mathbf{c}_1^T \cdot \text{db}^T) \in \mathbb{Z}_q^{\kappa \times \ell} \\ \mathbf{h} \leftarrow \text{ans}_1 \cdot \mathbf{A}_2 \in \mathbb{Z}_q^{\kappa \times n}, \\ \begin{bmatrix} \text{ans}_h \\ \text{ans}_2 \end{bmatrix} \leftarrow \begin{bmatrix} \text{hint}_s \\ \text{ans}_1 \end{bmatrix} \cdot \mathbf{c}_2 \in \mathbb{Z}_q^{\kappa(n+1)}. \end{cases}$
- Return  $\text{ans} \leftarrow (\mathbf{h}, \text{ans}_h, \text{ans}_2)$ .

Recover( $\text{st}, \text{hint}_c, \text{ans}$ )  $\rightarrow$   $d$ .

- Parse:
  - $(\mathbf{s}_1 \in \mathbb{Z}_q^n, \mathbf{s}_2 \in \mathbb{Z}_q^n) \leftarrow \text{st}$ ,
  - $(\mathbf{h} \in \mathbb{Z}_q^{\kappa \times n}, \text{ans}_h \in \mathbb{Z}_q^{\kappa \times n}, \text{ans}_2 \in \mathbb{Z}_q^\ell) \leftarrow \text{ans}$ .
- Compute:

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} &\leftarrow \begin{bmatrix} \text{ans}_h \\ \text{ans}_2 \end{bmatrix} - \begin{bmatrix} \text{hint}_c \\ \mathbf{h} \end{bmatrix} \cdot \mathbf{s}_2 \in \mathbb{Z}_q^{\kappa \times (n+1)}, \\ \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{a}_1 \end{bmatrix} &\leftarrow \text{Recomp} \left( \text{Round}_{\Delta} \left( \begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} \right) / \Delta \right) \in \mathbb{Z}_q^{n+1}, \\ \hat{d} &\leftarrow \mathbf{a}_1 - \mathbf{s}_1^T \mathbf{h}_1 \in \mathbb{Z}_q. \end{aligned}$$

- Return  $d \leftarrow \text{Round}_{\Delta}(\hat{d}) / \Delta \in \mathbb{Z}_p$ .

Figure 14: The DoublePIR protocol.

## E.2 Correctness and security of DoublePIR

We now prove the following theorem:

**Theorem E.1** (DoublePIR). *On database size  $N \in \mathbb{N}$ , correctness failure probability  $\delta$ , adversary runtime  $T$ , and security failure probability  $\epsilon$ , let*

- $\ell, m \in \mathbb{N}$  be database dimensions (such that  $\ell \cdot m \geq N$ ),
- $\chi$  be the discrete Gaussian distribution with variance  $\sigma^2$ ,
- $(n, q, \chi)$  be LWE parameters achieving  $(T, \epsilon)$ -security for  $\max(\ell, m)$  LWE samples, and
- $p \in \mathbb{N}$  be a plaintext modulus chosen to satisfy

$$\lfloor q/p \rfloor \geq \sigma p \sqrt{2 \max(m, \ell) \cdot \ln \left( \frac{2(\kappa(n+1)+1)}{\delta} \right)}, \quad (3)$$

where  $\kappa \in \mathbb{Z}$  is the scalar  $\left\lceil \frac{\log q}{\log p} \right\rceil$ .

Then, for random LWE matrices  $\mathbf{A}_1 \stackrel{\mathbb{R}}{\leftarrow} \mathbb{Z}_q^{m \times n}$  and  $\mathbf{A}_2 \stackrel{\mathbb{R}}{\leftarrow} \mathbb{Z}_q^{\ell \times n}$ , DoublePIR is a  $(T - O(n \cdot \ell + m), 4\epsilon)$ -secure PIR scheme on database size  $N$ , over plaintext space  $\mathbb{Z}_p$ , with correctness error  $\delta$ .

*Proof.* We prove correctness and security separately.

**Correctness.** Consider a client that queries for the database value at index  $(i_{\text{row}}, i_{\text{col}}) \in [\ell] \times [m]$ . We observe that the Recover routine in DoublePIR essentially runs  $\kappa \cdot (n+1) + 1$  instances of SimplePIR’s recovery routine, to recover  $\kappa \cdot (n+1) + 1$  elements in  $\mathbb{Z}_p$ , namely:

- one  $\mathbb{Z}_p$ -element that encodes the database value at index  $(i_{\text{row}}, i_{\text{col}})$  (in the “first level of PIR”),
- $\kappa n$  additional  $\mathbb{Z}_p$ -elements that encode the  $n$   $\mathbb{Z}_q$ -elements

in row  $i_{\text{row}}$  of the first-level hint matrix,  $\mathbf{A}_1^T \cdot \text{db}^T$  (in the “second level of PIR”), and

- $\kappa$  additional  $\mathbb{Z}_p$ -elements that encode the single  $\mathbb{Z}_q$ -element at position  $i_{\text{row}}$  in the first level answer vector,  $\mathbf{c}_1^T \cdot \text{db}^T$  (in the “second level of PIR”).

For the recovery to succeed, all  $\kappa \cdot (n+1) + 1$  invocations of SimplePIR’s recovery must succeed. To bound this probability, we compute the probabilities of two separate events:

- $E_1$ , the event that the recovery of the  $\kappa \cdot (n+1)$  elements in the “second level of PIR” fails, i.e., that

$$\begin{bmatrix} \mathbf{h}_1 \\ \mathbf{a}_1 \end{bmatrix} \neq \text{Recomp} \left( \begin{bmatrix} \text{hint}_s \\ \text{ans}_1 \end{bmatrix} \right) \cdot \mathbf{u}_{i_{\text{row}}},$$

where right-multiplication by  $\mathbf{u}_{i_{\text{row}}}$  extracts the  $i_{\text{row}}$ -th column of the matrix.

- $E_2$ , the event that the recovery of the single element in the “first level of PIR” fails, i.e., that

$$\text{Round}_\Delta(\hat{d})/\Delta \neq \text{db}[i_{\text{row}}, i_{\text{col}}].$$

We make use of the following two propositions:

*Proposition E.2: Under the parameters in Theorem E.1,*

$$\Pr[E_1] \leq \frac{\kappa(n+1)}{\kappa(n+1)+1} \cdot \delta.$$

*Proof.* DoublePIR’s recovery routine computes:

$$\begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} \leftarrow \text{Recomp} \left( \text{Round}_\Delta \left( \begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} \right) / \Delta \right) \in \mathbb{Z}_q^{n+1}.$$

Define  $\text{db}_2 = \begin{bmatrix} \text{hint}_s \\ \text{ans}_1 \end{bmatrix} \in \mathbb{Z}_q^{\kappa(n+1)+\ell}$  and let  $E$  be the event that

$$\text{Round}_\Delta \left( \begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} \right) / \Delta \neq \text{db}_2 \cdot \mathbf{u}_{i_{\text{row}}}. \quad (4)$$

We have that  $\Pr[E] \geq \Pr[E_1]$  as  $E_1 \subset E$ . By rewriting terms:

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{h}}_1 \\ \hat{\mathbf{a}}_1 \end{bmatrix} &= \begin{bmatrix} \text{ans}_h - \text{hint}_c \cdot \mathbf{s}_2 \\ \text{ans}_2 - \mathbf{h} \cdot \mathbf{s}_2 \end{bmatrix} \\ &= \begin{bmatrix} \text{hint}_s \cdot \mathbf{c}_2 - (\text{hint}_s \cdot \mathbf{A}_2) \cdot \mathbf{s}_2 \\ \text{ans}_1 \cdot \mathbf{c}_2 - (\text{ans}_1 \cdot \mathbf{A}_2) \cdot \mathbf{s}_2 \end{bmatrix} \\ &= \begin{bmatrix} \text{hint}_s \cdot (\mathbf{A}_2 \cdot \mathbf{s}_2 + \mathbf{e}_2 + \Delta \cdot \mathbf{u}_{i_{\text{row}}}) - \text{hint}_s \cdot \mathbf{A}_2 \cdot \mathbf{s}_2 \\ \text{ans}_1 \cdot (\mathbf{A}_2 \cdot \mathbf{s}_2 + \mathbf{e}_2 + \Delta \cdot \mathbf{u}_{i_{\text{row}}}) - \text{ans}_1 \cdot \mathbf{A}_2 \cdot \mathbf{s}_2 \end{bmatrix} \\ &= \text{db}_2 \cdot \mathbf{e}_2 + \Delta \cdot \text{db}_2 \cdot \mathbf{u}_{i_{\text{row}}}. \end{aligned}$$

Comparing this with Equation (4), we see that  $E$  happens if and only if  $\|\text{db}_2 \cdot \mathbf{e}_2\|_\infty \geq \Delta/2$ . The vector  $\mathbf{e} := \text{db}_2 \cdot \mathbf{e}_2$  consists of  $\kappa \cdot (n+1)$  elements, which we denote as  $e_1, \dots, e_{\kappa \cdot (n+1)}$ . As in our correctness proof for SimplePIR, we use a concentration inequality [68, Lemma 2.2][12, Lemma 2.4] on each of these terms: since  $\chi$  is the discrete Gaussian distribution with

variance  $\sigma^2 = \frac{s^2}{2\pi}$  for some  $s > 0$ , for every  $j \in [\kappa \cdot (n+1)]$  and for any  $T > 0$ , it holds that,

$$\Pr[|e_j| \geq T \cdot s \cdot \|\text{db}_2[j, :]\|] < 2 \exp(-\pi \cdot T^2),$$

where  $\|\cdot\|$  is the Euclidean norm. Note that since we can represent the elements in  $\text{db}_2$  in the range  $(-p/2, p/2]$ , we have that  $\|\text{db}_2[j, :]\| \leq \sqrt{\ell} \cdot (p/2) = \sqrt{\ell} \cdot p/2$ . Therefore, for every  $j \in [\kappa \cdot (n+1)]$ , for any  $T > 0$ , it holds that,

$$\Pr[|e_j| \geq T \cdot s \cdot \sqrt{\ell} \cdot p/2] < 2 \exp(-\pi \cdot T^2).$$

Now, we take  $T = \Delta/(2s \cdot \sqrt{\ell} \cdot p/2)$ , and we see that, for every  $j \in [\kappa \cdot (n+1)]$ ,

$$\begin{aligned} \Pr[|e_j| \geq \Delta/2] &< \frac{\delta}{\kappa(n+1)+1}, \text{ as long as} \\ 2 \exp\left(-\pi \cdot \left(\frac{\Delta}{2s \cdot \sqrt{\ell} \cdot p/2}\right)^2\right) &\leq \frac{\delta}{\kappa(n+1)+1}. \end{aligned}$$

Substituting  $\Delta = \lfloor q/p \rfloor$  and  $s = \sigma \cdot \sqrt{2\pi}$ , we rewrite the second equation as:

$$\lfloor q/p \rfloor \geq \sigma \cdot p \cdot \sqrt{2\ell \cdot \ln\left(\frac{2\kappa(n+1)+1}{\delta}\right)}.$$

By Equation (3), this condition holds. Thus, we have shown that, for each  $j \in [\kappa \cdot (n+1)]$ ,  $|e_j| < \Delta/2$  holds except with probability  $\frac{\delta}{\kappa(n+1)+1}$ . Therefore, by a union bound, all  $e_1, \dots, e_{\kappa \cdot (n+1)}$  have an absolute value smaller than  $\Delta/2$  with probability  $\frac{\kappa(n+1)}{\kappa(n+1)+1} \cdot \delta$ , which completes the proof.  $\square$

*Proposition E.3: Under the parameters in Theorem E.1,*

$$\Pr[E_2 | \neg E_1] \leq \frac{1}{\kappa(n+1)+1} \cdot \delta.$$

*Proof.* If  $\neg E_1$  occurs, we know that

$$\begin{aligned} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{a}_1 \end{bmatrix} &= \text{Recomp} \left( \begin{bmatrix} \text{hint}_s \\ \text{ans}_1 \end{bmatrix} \right) \cdot \mathbf{u}_{i_{\text{row}}} \\ &= \begin{bmatrix} \mathbf{A}_1^T \cdot \text{db}^T \\ \mathbf{c}_1^T \cdot \text{db}^T \end{bmatrix} \cdot \mathbf{u}_{i_{\text{row}}}. \end{aligned}$$

Therefore,

$$\begin{aligned} \hat{d} &= \mathbf{a}_1 - \mathbf{s}_1^T \cdot \mathbf{h}_1 \\ &= (\mathbf{c}_1^T \cdot \text{db}^T - \mathbf{s}_1^T \cdot \mathbf{A}_1^T \cdot \text{db}^T) \cdot \mathbf{u}_{i_{\text{row}}} \\ &= ((\mathbf{s}_1^T \cdot \mathbf{A}_1^T + \mathbf{e}_1^T + \Delta \mathbf{u}_{i_{\text{col}}}^T) \cdot \text{db}^T - \mathbf{s}_1^T \cdot \mathbf{A}_1^T \cdot \text{db}^T) \cdot \mathbf{u}_{i_{\text{row}}} \\ &= \mathbf{e}^T \cdot \text{db}^T \cdot \mathbf{u}_{i_{\text{row}}} + \Delta \mathbf{u}_{i_{\text{col}}}^T \cdot \text{db}^T \cdot \mathbf{u}_{i_{\text{row}}}. \end{aligned}$$

Thus, we see that Recover outputs the correct database element,  $\text{db}[i_{\text{row}}, i_{\text{col}}]$ , if and only if  $|\mathbf{e}^T \cdot \text{db}^T \cdot \mathbf{u}_{i_{\text{row}}}| < \Delta/2$ . By essentially the same analysis as in Proposition E.2, the event

that  $|\mathbf{e}^T \cdot \text{db}^T \cdot \mathbf{u}_{i_{\text{row}}}| \geq \Delta/2$  happens occurs with probability at most  $\frac{1}{\kappa(n+1)+1} \cdot \delta$ , given that the condition

$$\lfloor q/p \rfloor \geq \sigma \cdot p \cdot \sqrt{2m \cdot \ln \frac{2(\kappa(n+1)+1)}{\delta}}$$

holds, which is true by Equation (3).  $\square$

Combining the two propositions, we conclude that DoublePIR's correctness error is upper bounded by

$$\Pr[E_1 \vee E_2] \leq \Pr[E_1] + \Pr[E_2 | \neg E_1] \leq \delta.$$

**Security.** DoublePIR's security follows from the hardness of LWE with a reused LWE matrix [83]. We rely on the following lemma:

**Lemma E.4.** *Let  $N \in \mathbb{N}$  be the database size,  $\ell, m \in \mathbb{N}$  be the database dimensions (such that  $\ell \cdot m \geq N$ ),  $(n, q, \chi)$  be the LWE parameters, and  $\mathbf{A}_1 \in \mathbb{Z}_q^{m \times n}$  and  $\mathbf{A}_2 \in \mathbb{Z}_q^{\ell \times n}$  be the random LWE matrices used in DoublePIR. For any  $i \in [N]$ , we define the distribution*

$$Q_i = \{(\mathbf{A}_1, \mathbf{A}_2, \text{qu}_i) : \_, \text{qu}_i \leftarrow \text{Query}(i)\}.$$

*If the  $(n, q, \chi)$ -LWE problem with  $\max(\ell, m)$  samples is  $(T, \epsilon)$ -hard, then any algorithm running in time  $T - O(n \cdot \ell + m)$  has success probability at most  $2\epsilon$  in distinguishing  $Q_i$  from the distribution  $\mathcal{D} = \{(\mathbf{A}_1, \mathbf{A}_2, (\mathbf{r}_1, \mathbf{r}_2)) : \mathbf{r}_1 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^m, \mathbf{r}_2 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^\ell\}$ .*

*Proof.* Consider any index  $i \in [N]$ , and decompose  $i$  into the pair of coordinates  $(i_{\text{row}}, i_{\text{col}}) \in [\ell] \times [m]$  as done in DoublePIR's Query routine. Let  $\mathbf{u}_{i_{\text{col}}} \in \mathbb{Z}_q^m$  denote unit vector  $i_{\text{col}}$  in  $\mathbb{Z}_q^m$  and  $\mathbf{u}_{i_{\text{row}}} \in \mathbb{Z}_q^\ell$  denote unit vector  $i_{\text{row}}$  in  $\mathbb{Z}_q^\ell$ . We define the following distribution:

$$\mathcal{H}_i = \{(\mathbf{A}_1, \mathbf{A}_2, (\mathbf{r}_1, \mathbf{A}_2 \mathbf{s}_2 + \mathbf{e}_2)) : \mathbf{r}_1 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^m, \mathbf{s}_2 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^n, \mathbf{e}_2 \leftarrow^{\mathbb{R}} \chi^\ell\}$$

Now, consider any algorithm  $\mathcal{A}$  that runs in time  $t$  and distinguishes between distributions  $Q_i$  and  $\mathcal{D}$  with probability  $p$ . Let  $p_0, p_1$ , and  $p_2$  be the probabilities that  $\mathcal{A}$  outputs 1 when given samples from  $Q_i, \mathcal{H}_i$ , and  $\mathcal{D}$  respectively. Thus, by definition, it must hold that  $|p_0 - p_2| = p$ . From the triangle inequality we see that  $|p_0 - p_1| + |p_1 - p_2| \geq p$ .

Then, using  $\mathcal{A}$  as a subroutine, we build an algorithm  $\mathcal{B}_1$  that distinguishes between the  $(n, q, \chi)$ -LWE problem with  $m$  samples and the uniform distribution, in time  $t + O(m \cdot n + \ell)$  and with success probability  $|p_0 - p_1|$ . Because the  $(n, q, \chi)$ -LWE problem with  $m$  samples is  $(T, \epsilon)$ -hard, if  $t \leq T - O(m \cdot n + \ell)$ , it must hold that  $|p_0 - p_1| \leq \epsilon$ . Similarly, we build an algorithm  $\mathcal{B}_2$  that distinguishes between the  $(n, q, \chi)$ -LWE problem with  $\ell$  samples and the uniform distribution in time  $t + O(m + \ell)$  and with success probability  $|p_1 - p_2|$ . Because the  $(n, q, \chi)$ -LWE problem with  $\ell$  samples is  $(T, \epsilon)$ -hard, if  $t \leq T - O(m + \ell)$ , it must hold that  $|p_1 - p_2| \leq \epsilon$ .

Thus, if  $t \leq T - O(m \cdot n + \ell)$ , we see that  $2\epsilon \geq |p_0 - p_1| + |p_1 - p_2| \geq p$ . We have shown that any algorithm  $\mathcal{A}$  that runs

in time  $T - O(m \cdot n + \ell)$  can distinguish between  $Q_i$  and  $\mathcal{D}$  with probability at most  $2\epsilon$ .

**Algorithm 3**  $\mathcal{B}_1$ , on input  $\mathbf{A}_1, \mathbf{b}_1$ :

---

Compute  $\mathbf{c}_1 \leftarrow \mathbf{b}_1 + \lfloor q/p \rfloor \cdot \mathbf{u}_{i_{\text{col}}}$ .  
 Sample  $\mathbf{s}_2 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^n, \mathbf{e}_2 \leftarrow^{\mathbb{R}} \chi^\ell$ .  
 Compute  $\mathbf{c}_2 \leftarrow (\mathbf{A}_2 \mathbf{s}_2 + \mathbf{e}_2) + \lfloor q/p \rfloor \cdot \mathbf{u}_{i_{\text{row}}}$ .  
 Output  $\mathcal{A}(\mathbf{A}_1, \mathbf{A}_2, (\mathbf{c}_1, \mathbf{c}_2))$ .

---

**Algorithm 4**  $\mathcal{B}_2$ , on input  $\mathbf{A}_2, \mathbf{b}_2$ :

---

Sample  $\mathbf{r}_1 \leftarrow^{\mathbb{R}} \mathbb{Z}_q^m$ .  
 Compute  $\mathbf{c}_2 \leftarrow \mathbf{b}_2 + \lfloor q/p \rfloor \cdot \mathbf{u}_{i_{\text{row}}}$ .  
 Output  $\mathcal{A}(\mathbf{A}_1, \mathbf{A}_2, (\mathbf{r}_1, \mathbf{c}_2))$ .

---

$\square$

Lemma E.4 states that any algorithm running in time  $T - O(n \cdot \ell + m)$  can distinguish the queries made by a client in DoublePIR from random vectors with success probability at most  $2\epsilon$ . Therefore, any algorithm running in time  $T - O(n \cdot \ell + m)$  can distinguish the queries made by a client in DoublePIR to any pair of indices  $i \in [N]$  and  $j \in [N]$  with success probability at most  $4\epsilon$ . Thus, we have shown that DoublePIR is  $(T - O(n \cdot \ell + m), 4\epsilon)$ -secure.  $\square$

DoublePIR also meets the  $Q$ -query security definition that we give in Appendix C.2. A hybrid argument shows that:

**Corollary E.5** ( $Q$ -query security of DoublePIR). *For any sequence  $I \in [N]^Q$  of  $Q$  indices, we define  $\mathcal{D}_I$  to be the query distribution that it induces, i.e.,*

$$\mathcal{D}_I = \{(\mathbf{A}_1, \mathbf{A}_2, \text{qu}_k) : \_, \text{qu}_k \leftarrow \text{Query}(I_k)\}_{k \in [Q]}$$

*If the  $(n, q, \chi)$ -LWE problem with  $\max(\ell, m)$  samples is  $(T, \epsilon)$ -hard, then  $\mathcal{D}_I$  is  $(T - O(Q \cdot (n \cdot \ell + m)), 2Q\epsilon)$ -indistinguishable from the random query distribution.*

### E.3 Additional extensions and optimizations

The SimplePIR optimizations involving fast matrix multiplication algorithms (that we outline in Appendix C.3) apply to DoublePIR as well. We list some further extensions and efficiency improvements to DoublePIR in this section.

**Handling database updates.** When the database contents change, the server in DoublePIR needs to update its preprocessed hints,  $\text{hint}_s$  and  $\text{hint}_c$  (see Remark 3.1). In more detail, when  $c$  rows of the database matrix change (for some number  $c \geq 0$ ), the server needs to:

- perform  $O(c \cdot m \cdot n + c \cdot \kappa \cdot n)$  operations in  $\mathbb{Z}_q$ , and
- send  $\min(c \cdot \kappa \cdot n, \kappa \cdot n^2)$  elements in  $\mathbb{Z}_q$  back to the client.

This update procedure works as follows. In DoublePIR, the server computes the hints as  $\text{hint}_s \leftarrow \text{Decomp}(\mathbf{A}_1^T \cdot \text{db}^T)$  and  $\text{hint}_c \leftarrow \text{hint}_s \cdot \mathbf{A}_2$ . For each row  $\mathbf{d} \in \mathbb{Z}_q^m$  of the database that has changed, the server can update the corresponding columns of  $\text{hint}_s$  to be  $\mathbf{u} \leftarrow \text{Decomp}(\mathbf{A}_1^T \cdot \mathbf{d}^T) \in \mathbb{Z}_q^{\kappa n}$ . To update  $\text{hint}_c$ ,

for each row of the database that has changed, the server would then have to add the outer product between the column vector,  $\mathbf{u}$ , and the corresponding row of  $\mathbf{A}_2$ , to  $\text{hint}_c$ . The server can perform this computation and send the updated  $\text{hint}_c$  back to the client; or, if fewer than  $n$  rows of the database were changed, the server can send only the associated column vectors,  $\mathbf{u}$ , back to the client (to save on bandwidth) and have the client update its  $\text{hint}_c$  locally. Again, additions and deletions of rows from the database can be handled similarly.

**Reducing the client storage, when the client makes a bounded number of queries.** If the client knows ahead of time that it will only make  $Q \ll n$  queries, it can reduce its local storage to  $Q \cdot \kappa \cdot n$  (rather than  $\kappa \cdot n^2$ ) elements in  $\mathbb{Z}_q$ . To do so, the client samples  $Q$  pairs of LWE secrets  $(s_1, s_2)$  in advance (e.g., using a pseudorandom function). Then, rather than store  $\text{hint}_c$  in its entirety, the client precomputes and stores only the information that it will need for recovery: for each of the  $Q$  queries it intends to make, the client precomputes  $\text{hint}_c \cdot s_2 \in \mathbb{Z}_q^{\kappa n}$  (where  $s_2$  denotes the secret key used for the corresponding query). Finally, the client discards  $\text{hint}_c$ .

In theory, the client storage could be reduced even further by applying the local rounding trick (see Appendix C.3) to the precomputed information,  $\text{hint}_c \cdot s_2$ , stored by the client. After doing so, the client would store only a single element in  $\mathbb{Z}_q$ —namely,  $s_1^T \cdot \text{Recomp}(\text{Round}_\Delta(\text{hint}_c \cdot s_2))$ —for each query it intends to make. However, doing so would require using a much larger value of  $\lfloor q/p \rfloor$  to maintain correctness and is thus not profitable in our parameter regime.

## F Additional material on our private set-membership data structure

### F.1 Syntax and properties

We define the syntax and security properties for our approximate private set membership data structure. The data structure is parameterized by a universe  $\mathcal{U} \subseteq \{0, 1\}^*$  of possible strings and consists of two routines:

$\text{Setup}(S) \rightarrow D$ . Take as input a set of strings  $S \subseteq \mathcal{U}$  and output a data structure  $D \in \{0, 1\}^{a \times b}$ , where  $a, b \in \mathbb{N}$  can depend on  $|S|$ .

$\text{Query}(\sigma) \rightarrow (i, j)$ . Given a candidate string  $\sigma \in \mathcal{U}$ , output an index  $(i, j) \in [a] \times [b]$  in the data structure  $D$  to probe. If  $\sigma \in S$ , it holds that  $D_{i,j} = 1$ .

**Properties.** For a set  $S \subseteq \mathcal{U}$  and string  $\sigma \in \mathcal{U}$ , define

$$\text{Accept}(S, \sigma) := \Pr \left[ D_{i,j} = 1 : \begin{array}{l} D \leftarrow \text{Setup}(S) \\ (i, j) \leftarrow \text{Query}(\sigma) \end{array} \right]$$

where the probability is taken *only* over the randomness of the Query algorithm.

**Correctness.** An approximate membership test is *correct* if, for all  $S \subseteq \mathcal{U}$  and  $\sigma \in S$ ,  $\text{Accept}(S, \sigma) = 1$ .

**Constant adversarial false-positive rate.** An approximate membership test has *adversarial false-positive rate*  $\epsilon$  if, for all  $S \subseteq \mathcal{U}$  and  $\sigma \notin S$ ,  $\text{Accept}(S, \sigma) \leq \epsilon$ .

**Oblivious reads.** An approximate membership test supports *oblivious reads* if the row of the data structure that the Query algorithm probes reveals no information about the query string  $\sigma$ . In other words, for all strings  $\sigma_0, \sigma_1 \in \mathcal{U}$ , the following distributions of the row indices are identical:

$$\{i : (i, \_) \leftarrow \text{Query}(\sigma_0)\} \equiv \{i : (i, \_) \leftarrow \text{Query}(\sigma_1)\}.$$

In a deployment, to probe the element at row  $i$  and column  $j$  in the data structure, the client sends  $i$  in the clear, along with a PIR query for  $j$ , to the server. The server executes the PIR scheme over the  $i$ -th row, and sends the PIR answer (and, if needed, the client hint for that row) back to the client. As the reads are oblivious, the server learns nothing about  $\sigma$ . Performing PIR over only a row of the data structure (rather than the whole data structure) will be the source of our  $\lambda$ -fold performance improvements compared to related approaches.

### F.2 Our approximate membership test

Our construction is parameterized by integers  $a, k \in \mathbb{N}$  and a set size  $N$ . Our data structure then maps the set of strings  $S$  into a matrix of  $a$ -by- $(kN)$  bits. The construction uses cryptographic hash functions  $H_1, \dots, H_a : \mathcal{U} \rightarrow [kN]$ , which we model as independent random oracles [14] and which are chosen *after* the set  $S$  has been fixed.

$\text{Setup}(S) \rightarrow D \in \{0, 1\}^{a \times kN}$ :

- Let  $D \in \{0, 1\}^{a \times kN}$  be the all-zeros matrix.
- For each  $i \in [a]$  and  $\sigma \in S$ , set  $D_{i, H_i(\sigma)} = 1$ .
- Return  $D$ .

$\text{Query}(\sigma) \rightarrow (i, j) \in [a] \times [kN]$ :

- Choose  $i \xleftarrow{\mathbb{R}} [a]$ , and output  $(i, H_i(\sigma))$ .

### F.3 Security analysis

*Proposition F.1:* The approximate membership-test data structure (Setup, Query), on parameters  $a \geq 2(\log(|\mathcal{U}|) + \lambda)$  and  $k \geq 8$ , is correct, has false-positive rate  $1/2$ , and has oblivious reads. The construction fails with probability  $2^{-\lambda}$ , over the choice of the hash functions modeled as independent random oracles. Concretely, on  $|\mathcal{U}| = 2^{256}$ , taking  $k = 8$  and  $a = 768$  gives false-positive rate  $1/2$  and failure probability  $2^{-128}$ .

*Proof.* Correctness and oblivious reads follow by construction. Constant adversarial false-positive rate is proved in Lemma F.2, which follows.  $\square$

**Lemma F.2.** The approximate membership-test data structure (Setup, Query) on parameters  $a$  and  $k$  has adversarial false-positive rate at most  $c$  and fails with probability  $2^{-\lambda}$ , for any  $c \in (0, 1)$  that satisfies  $|\mathcal{U}| \binom{a}{c} k^{-ca} \leq 2^{-\lambda}$ .

When  $k \geq 8$  and taking  $c = 1/2$ , we get a data structure with adversarial false-positive rate  $1/2$  and failure probability at most  $2^{-\lambda}$ , as long as  $|\mathcal{U}| 2^{-a/2} \leq 2^{-\lambda}$ .

*Proof.* Fix any set  $S \subseteq \mathcal{U}$  and consider  $D \leftarrow \text{Setup}(S)$ . Consider any  $c \in (0, 1)$  that satisfies  $|\mathcal{U}| \binom{a}{ca} k^{-ca} \leq 2^{-\lambda}$ .

For string  $\sigma \in \mathcal{U} \setminus S$ , let  $\text{Bad}_\sigma$  be the event, over the random choice of the random oracle, that there is a set of rows  $R \subseteq [a]$  of size  $ca$ , such that for all  $i \in R$ ,  $D_{i, H_i(\sigma)} = 1$ . First, we bound the probability, over the choice of the random oracle, that there exists any bad string  $\sigma$ .

In particular, fix a string  $\sigma \in \mathcal{U} \setminus S$ . Now we analyze the event  $\text{Bad}_\sigma$  as follows. Each row of  $D$  contains at most  $|S| = N$  ones and has  $kN$  entries total. Thus, an independently sampled random element in  $D$  is non-zero with probability at most  $k^{-1}$ . For a particular set  $R \subseteq [a]$  of size  $ca$ , let  $\text{Bad}_{\sigma, R}$  be the event that for all  $i \in R$ ,  $D_{i, H_i(\sigma)} = 1$ . Since the hash functions are modelled as random oracles, we have that  $\Pr[\text{Bad}_{\sigma, R}] \leq k^{-ca}$ . By a union bound over all possible sets  $R$ ,

$$\Pr[\text{Bad}_\sigma] = \Pr \left[ \bigvee_R \text{Bad}_{\sigma, R} \right] \leq \binom{a}{ca} \cdot k^{-ca}$$

Let  $\text{Bad} := \bigvee_{\sigma \in \mathcal{U} \setminus S} \text{Bad}_\sigma$ . By a union bound over all possible strings  $\sigma$ , we get that:

$$\Pr[\text{Bad}] \leq |\mathcal{U}| \cdot \Pr[\text{Bad}_\sigma] \leq |\mathcal{U}| \cdot \binom{a}{ca} \cdot k^{-ca}.$$

So, with probability  $1 - |\mathcal{U}| \binom{a}{ca} k^{-ca}$  over the choice of the random oracles, there do not exist any bad strings  $\sigma$ .

Now, consider the event  $\neg \text{Bad}$  (i.e., there exists no bad string). In this case, for any  $\sigma \in \mathcal{U} \setminus S$ , we know that, with probability at most  $c$  (over the random coins of Query), the query algorithm will sample a row  $i$  such that  $D_{i, H_i(\sigma)} = 1$ . Thus, we know that:

$$\text{Accept}(S, \sigma) = \Pr_i [D_{i, H_i(\sigma)} = 1] \leq c,$$

where the probability is taken over the random choice of  $i$  by the Query algorithm.

Finally, taking  $k \geq 8$  and  $c = 1/2$ , we see that  $\binom{a}{a/2} \leq 2^a$  and  $k^{-ca} \leq 2^{-3a/2}$ , which completes the proof.  $\square$

## G Additional implementation material

In Figure 15, we give a code listing for the core of SimplePIR (without the extensions of Section 4.3). We do not include the code for the following helper functions:

- `MatrixRand`, which takes as input a matrix height, a matrix width, and the modulus to be used for the matrix entries, and returns a uniformly random matrix of the given dimensions.
- `MatrixGaussian`, which takes as input a matrix height, a matrix width, and the modulus to be used for the matrix

```
// Data types
type Matrix struct {
    rows uint
    cols uint
    data []uint
}

type Params struct {
    n          uint // LWE secret dimension
    db_height  uint // DB height
    db_width   uint // DB width
    q          uint // ciphertext modulus
    p          uint // plaintext modulus
}

// SimplePIR methods
func Init(p Params) Matrix {
    A := MatrixRand(p.db_height, p.n, p.q)
    return A // 'A' is a public parameter
}

func Setup(DB Matrix, A Matrix) Matrix {
    H := MatrixMul(DB, A)
    return H // 'H' is the client hint
}

func Query(i uint, A Matrix,
           p Params) (Matrix, Matrix) {
    s := MatrixRand(p.n, 1, p.q)
    err := MatrixGaussian(p.db_width, 1, p.q)
    query := MatrixMul(A, s)
    query.MatrixAdd(err)
    query[i % p.db_width] += (p.q / p.p)
    return s, query // 's' is the client state
                    // 'query' is the client query
}

func Answer(DB Matrix, query Matrix) Matrix {
    ans := MatrixMul(DB, query)
    return ans // 'ans' is the server's answer
}

func Recover(i uint64, H, ans, s Matrix,
            p Params) uint64 {
    interm := MatrixMul(H, s)
    ans.MatrixSub(interm)
    noised := ans.data[i / p.db_width]
    denoised := Round(noised, (p.q / p.p))
    return denoised // the recovered database value
}
```

Figure 15: SimplePIR Go code.

entries, and returns a matrix of the given dimensions whose entries were sampled from the discrete gaussian distribution (with the appropriate, hard-coded standard deviation).

- `MatrixMul`, which takes as input two matrices and returns their product, over the appropriate field.
- `MatrixAdd`, which takes as input two matrices and returns their sum, over the appropriate field.
- `MatrixSub`, which takes as input two matrices and returns their difference, over the appropriate field.
- `Round`, which takes as input a value and a base and returns the value rounded to the nearest multiple of the base, and then divided by the base.