

# FrodoPIR: Simple, Scalable, Single-Server Private Information Retrieval

Alex Davidson, Gonçalo Pestana, Sofia Celi

Brave Software

**Abstract.** We design FrodoPIR — a highly configurable, *stateful*, single-server Private Information Retrieval (PIR) scheme that involves an offline phase that is completely *client-independent*. Coupled with small online overheads, it leads to much smaller amortized financial costs on the server-side than previous approaches. In terms of performance for a database of 1 million 1KB elements, FrodoPIR requires  $< 1$  second for responding to a client query, has a server response size blow-up factor of  $< 3.6\times$ , and financial costs are  $\sim \$1$  for answering 100,000 client queries. Our experimental analysis is built upon a simple, non-optimized Rust implementation, illustrating that FrodoPIR is particularly suitable for deployments that involve large numbers of clients.

## 1 Introduction

A Private Information Retrieval (PIR) scheme provides the ability for clients to retrieve items from an online database, without revealing anything about their queries to the untrusted host server(s). Applications of practical PIR schemes include: anonymous communication [7, 61], anonymous media streaming [47], privacy-preserving ad-delivery [45, 68, 63], private location discovery [37], private contact discovery [16], password-checking [3], and SafeBrowsing [54]. PIR schemes are split into those that are information-theoretically secure, but require the database to be shared between multiple non-colluding servers [5, 26, 9, 11, 31, 10, 12, 75, 36, 35, 28, 40, 72, 58]; and those that are computationally-secure against a single untrusted server [3, 31, 6, 21, 24, 38, 55, 56, 1, 64, 65, 62, 59].

Multi-server PIR constructions are typically more efficient than single-server schemes. However, finding non-colluding servers to jointly fulfill the PIR functionality can be unrealistic and burdensome. To avoid such problems, developing practical single-server PIR schemes is a desirable goal. The most efficient single-server PIR schemes are based on fully homomorphic encryption (FHE), with security derived from the ring learning with errors (RLWE) assumption [1, 6, 62, 3, 64, 59]. Unfortunately, these schemes incur computational, bandwidth, and consequent financial overheads for answering client queries on standard, cloud-based infrastructure that would make them expensive to run at scale. Even the most efficient require several seconds to process a single query on a database of 1 million 1KB elements.

To drive down online and financial costs, a recent line of work of single-server PIR moves large proportions of the expensive online computation and

communication to an offline phase [62, 65, 30] (a technique that also applies in the two-server model [31, 58]). In this model, the client and server prepare an offline internal state to be used for making online queries. Such schemes are referred to as *offline-online* or *stateful*, as opposed to *online-only* or *stateless*. Many works [62, 65, 30] have focused on developing PIR schemes with efficient online phases. The recent work of Corrigan-Gibbs et al. [30], for example, produces a stateful single-server PIR scheme with sublinear efficiency costs.

A key difficulty that has gone unsolved is that either the computation or the communication costs induced during the offline phase scale linearly in the number of clients that will make queries [62, 30, 65]. Moreover, previous schemes require each individual client to make large numbers of queries (e.g.  $\sqrt{m}$  for  $m$  DB elements) to ensure that the amortized costs remain sublinear. Ultimately, this still results in significant financial costs for any server that plans to run a PIR service in standard cloud-based infrastructure and that will answer queries from large numbers of clients. As a consequence, single-server PIR remains unusable in many real-world applications.

*Our results.* We build FrodoPIR: a stateful PIR scheme that is built directly upon the learning with errors (LWE) problem only, rather than using RLWE and FHE-based technologies. Similarly to FrodoKEM with respect to lattice-based key exchange [17], we show that — counter to accepted intuition — eschewing ring lattice structures can lead to flexible and practically efficient PIR schemes. The main benefit of FrodoPIR is that the offline phase of the protocol is performed by the server alone, completely independent of the number of clients or queries that will be made. This results in low amortized computation overheads, and an offline client download size that is a tiny fraction of the entire server database.

Our results highlight that the current bottleneck for deploying practical *stateful* PIR schemes is heavily related to the per-client scalability of the offline preprocessing phase. Previous schemes have optimized primarily for per-client asymptotics, which we show do not necessarily translate into financially cheap real-world systems. To this end, FrodoPIR represents an initial exploration in developing stateful PIR schemes that are suitable for large, real-world deployments, where lowering financial costs for server-side operators is of paramount importance. On top of this, FrodoPIR is significantly simpler than previous schemes, making no use of FHE techniques and requiring only modular arithmetic that can be implemented using standard 32-bit unsigned integer instructions. Our formal contributions are as follows:

1. A stateful single-server PIR scheme, known as FrodoPIR, with security derived from LWE.
2. A simple, open-source Rust implementation — containing only a few hundred lines of code.<sup>1</sup>
3. Experimental analysis that illustrates that FrodoPIR is cheaper to run in large multi-client deployments than all previous single-server PIR schemes.

---

<sup>1</sup> <https://github.com/brave-experiments/frodo-pir>

4. Detailed analysis of various configuration trade-offs and optimizations for FrodoPIR.

## 2 Background

### 2.1 Overview of Prior Approaches

PIR was first introduced as a cryptographic primitive by Chor, Gilboa, Kushilevitz, and Sudan [28]. Information-theoretic PIR (itPIR) sees the client interact with multiple non-colluding servers, each of which have access to some form of the same database, and the client combines the responses from each server locally [5, 26, 9, 11, 31, 10, 12, 75, 36, 35, 40, 72, 58]. Computationally-secure PIR (cPIR) relies only on a single-server, and provides computational security based on cryptographic assumptions [3, 31, 6, 21, 24, 38, 55, 56, 1, 64, 65, 62]. While itPIR schemes are more efficient, real-world systems that provide non-collusion guarantees prove very hard to devise in practice. Thus, we focus on cPIR henceforth.

*Stateless PIR.* Initial constructions of PIR schemes followed the framework of Kushilevitz and Ostrowsky [55], using additively homomorphic encryption (from number-theoretic assumptions) for hiding the client query [21, 24, 38, 56]. Such schemes are known as online-only or stateless, since the client does not have to store any information in order to launch queries. Stateless single-server PIR schemes of this nature have the following underlying structure.

- To learn the  $i^{\text{th}}$  DB element  $\text{DB}[i]$ , a client sends a vector  $\mathbf{v}$  of  $m$  additively homomorphic ciphertexts, where  $\mathbf{v}[i]$  encrypts 1 and all others encrypt 0.
- The server responds with a vector  $\mathbf{w}$ , where  $\mathbf{w}[j] = \mathbf{v}[j] * \text{DB}[j]$  ( $j \in [m]$ ,  $*$  denotes scalar multiplication and  $m$  denotes the number of DB elements).
- The client decrypts  $\mathbf{w}[i]$  and learns  $\text{DB}[i]$ .

Sion and Carbunar showed that such schemes actually perform much worse than simply having the client download the entire server database (DB), when the network bandwidth is just a few hundred Kbps [69]. This is a result of performing  $O(m)$  expensive arithmetic operations (modular exponentiations or multiplications) for every client query.

The results of [69] stood as a reference point for nearly a decade, until Aguilar-Melchor et al. [1] used lattice-based cryptography (inherently faster than number-theoretic approaches) to construct efficient single-server PIR. In their scheme, XPIR, the server computation time is approximately  $> 5$  seconds for a DB with  $m = 2^{20}$  elements, even with the aforementioned asymptotic overheads. Accordingly, bandwidth requirements for the client query are 18MB, and 590KB for the server response. Various schemes since have used RLWE-based FHE to propose similar schemes or optimizations of these methods, such as [33, 6, 64, 3, 62, 59]. In particular, the works of [6, 62, 3, 64, 59] exhibit various optimizations that transform the client query and server database to reduce the size of the query and server response (to around 64KB and 128KB, respectively), whilst maintaining similar or improving computational costs.

*Stateful PIR.* Unfortunately, stateless cPIR schemes still require computational overheads that are difficult to justify in a large-scale deployment. For example, to respond to a single client query for a database of 1 million 256B entries, it takes  $> 1$  second, and requires downloading at least tens of kilobytes of data [6, 59]. Such approaches are unlikely to scale for large numbers of clients, or in situations that require timely responses. Recent work has observed that online performance can be improved by moving expensive, query-independent computation to an initial offline phase [65, 62, 31, 30]. This allows reducing the online costs, as well as amortizing the costs of the offline phase across a number of client queries.

The scheme of Patel et al., known as PSIR [65], enjoys a very fast online phase, though this approach requires the client to download the entire server database in an offline phase — which violates fundamental PIR efficiency criterion: the total client communication remains smaller than downloading the entire database (Definition 5). The scheme of Mughees et al., known as Stateful OnionPIR (henceforth SONionPIR) [62], provides a (financially) cheaper approach than PSIR, but at the cost of large computational overheads during the offline phase, which is executed as a protocol between each client and the server. Thus, financial costs will scale linearly in the global number of client queries that are launched. While the single-server scheme of Corrigan-Gibbs and Kogan [31] has similar issues as SONionPIR, the very recent work of Corrigan-Gibbs et al. [30] constructs a PIR scheme (henceforth CHKPIR) where all (amortized) asymptotic complexities are sublinear in the number of DB elements  $m$ . Specifically, costs are  $O(\sqrt{m})$ , when clients make  $\sqrt{m}$  queries. Previous schemes require  $O(m)$  (symmetric) online operations. This reduces further the online costs, but the costs of the offline phase are very similar to the previous works of [65, 62]. In summary, the expensive offline phase in each scheme — that only amortizes per a single client’s queries — quickly becomes the main driver of the server-side costs.

The general idea behind each of [65, 62, 30] is that each client and the server cooperatively run a *private batch sum retrieval* protocol that samples  $c$  random subsets  $S_1, \dots, S_c$  of elements DB, and computes the sum  $s_i$  of all of the elements in each  $S_i$  and provides it to the client. During the online phase, the client that wants to query for the element  $e_j = \text{DB}[j]$  picks the first  $t \in [c]$ , where  $e_j \notin S_t$ . They then construct a partition  $\mathcal{P} = (P_1, \dots, P_k)$  of the indices of DB, where  $P_j = S$ , and send a succinct description of this partition to the server. The server expands each partition into the set of sums  $s_{P_1}, \dots, s_{P_k}$ . The client uses an underlying single-server PIR scheme to learn the sum  $s_{P_j}$ , and, finally, outputs  $s_{P_j} - s_t$  to learn  $e_j$ .

The PSIR scheme implements the private batch sum retrieval protocol by streaming the entire database to the client, while the SONionPIR and CHKPIR schemes both involve the client specifying their random subsets as FHE ciphertexts, and having the server construct each of the sums using homomorphic properties. When instantiating the underlying single-server PIR scheme during the online phase using FHE-based schemes (such as SealPIR, or *stateless* OnionPIR), it has been shown that stateful PIR result in much more efficient online phases

and significantly smaller server costs, when compared with stateless, online-only schemes [65, 62].

*Other privacy-preserving data access primitives.* Oblivious RAM (ORAM) provides data access pattern privacy for client queries to a server database [42, 43]. This problem is related to PIR, but provides privacy also for the server database: the client learns the queried DB element and nothing more. While recent ORAM schemes enjoy sublinear computation and communication [25, 32, 67, 70], none are inherently multi-client and this leads to very expensive real-world overheads.

Hamlin et al. present Private Anonymous Data Access (PANDA) [48], based on a symmetric-key formulation of PIR known as doubly-efficient PIR [13, 19, 23]. Doubly-efficient PIR schemes are similar to stateful schemes, where there is an initial phase that preprocesses the server database, but the online phase is totally stateless. Unfortunately, symmetric-key doubly-efficient PIR is inherently not multi-client. Public-key instantiations use multiple-servers [13], or are based on expensive cryptographic obfuscation [19]. Batch PIR [52, 49, 8, 57] uses batch codes to achieve sublinear amortized efficiency, by allowing clients to retrieve multiple items at once. Unfortunately, such savings do not apply in settings where queries are made *adaptively* — i.e. based on the results of previous queries — which is required in most standard applications. As such, we focus on developing efficient PIR schemes for handling adaptive client queries.

## 2.2 Limitations of Existing Stateful PIR Schemes

*Expensive preprocessing.* The key limitation of SOnionPIR and CHKPIR is the computational cost of the private batch sum retrieval protocol that takes place during the offline phase. This protocol must be invoked per-client, and involves at least  $O(m)$  server-side operations and  $O(\sqrt{m})$  communication ( $m = |\text{DB}|$ ). These costs are amortized across the number of queries  $c$  launched but, even after amortization, the computational costs remain large. For a DB of  $2^{20}$  1KB elements, the offline phase of SOnionPIR takes 25 seconds *per client query*.<sup>2</sup> For large multi-client systems, the potential for amortization diminishes and these costs quickly become prohibitive.

In contrast, in PSIR clients simply download the entire server database before only storing  $O(c)$  data; this results in multiple issues. First, as shown in [62], for large numbers of clients, the download cost becomes prohibitively large from a financial perspective, and will continue growing for larger databases and items [62]. Second, the PSIR approach is unable to satisfy the fundamental efficiency criterion required of PIR schemes (Definition 5).

*High online bandwidth consumption.* As a result of using FHE-based single-server PIR during the online phase, both SOnionPIR and PSIR have online server response sizes that are relatively very large compared to the size of the

<sup>2</sup> While CHKPIR has not been implemented, the offline phase is very similar and thus will incur similarly large computational overheads.

**Table 1.** Asymptotic comparison focusing on the dependency on the number of database elements,  $m$ , of practical approaches for realizing single-server PIR with adaptive queries (i.e. not including batch PIR schemes, logarithmic factors are ignored). All costs are amortized according to  $C$  clients that launch  $c = \sqrt{m}$  queries ( $m = |\text{DB}|$  is the total number of elements in the server database). Communication costs relate to the amount of data that is *sent* to the party. The financial costs are given relative to a database containing  $2^{20}$  1KB elements, are amortized per-query and per-client, and are calculated assuming a server that operates the same AWS EC2 architecture specified in Section 6. <sup>†</sup>The costs of CHKPIR are assumed to be zero for the online phase, and are thus completely dominated by the offline phase, which can be implemented using techniques from [65, 62, 30].

Approach	Security assumptions	Client costs					Server costs				Financial
		Communication		Computation		Storage	Communication		Computation		
		Offline	Online	Offline	Online		Offline	Online	Offline	Online	
<b>Stateless</b> [6, 62, 3]	RLWE	—	$m$	—	$m$	—	—	1	—	$m$	$\$5.2 \times 10^{-3}$
<b>PSIR</b> [65]	RLWE	—	1	$m$	$\sqrt{m}$	$\sqrt{m}$	$ \text{DB} /\sqrt{m}$	1	—	$m$	$\$8.8 \times 10^{-5}$
<b>SONionPIR</b> [62]	RLWE	$\sqrt{m}$	1	$k \cdot \sqrt{m}$	$k$	$\sqrt{m}$	$\sqrt{m}$	1	$\sqrt{m}$	$m$	$\$6.4 \times 10^{-4}$
<b>CHKPIR</b> [30]	RLWE	$\sqrt{m}$	$\sqrt{m}$	$\sqrt{m}$	$\sqrt{m}$	$\sqrt{m}$	$\sqrt{m}$	1	$\sqrt{m}$	$\sqrt{m}$	$\sim \$8.8 \times 10^{-5\dagger}$
<b>FrodoPIR</b>	LWE	—	$m$	$m$	1	$\lambda$	$\lambda \cdot m^{-1/2}$	1	$\sqrt{m}/C$	$m$	$\$(1.9/C \times 10^{-2} + 1.3 \times 10^{-5})$

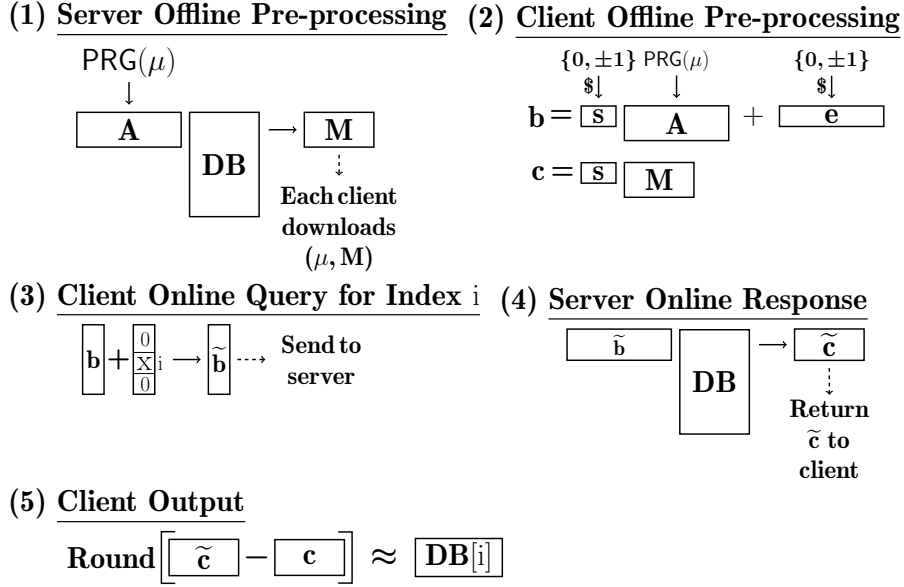
queried DB element. For example, for 1KB data elements, the response blow-up in SONionPIR is  $128\times$ , while in PSIR it is  $320\times$ . The work of CHKPIR uses similar underlying primitives and thus results in similar communication overheads.

*Practical security parameters.* PSIR and SONionPIR provide 115 and 111 bits of security, respectively [62, 65] using the primal-USVP cost model for estimating the hardness of cryptographic lattices, as shown in [2]. Achieving 128 bits of security can be important in cases where cryptographic tools must satisfy regulatory compliance checks. Increasing the concrete security parameters of either scheme would require modifying the LWE parameters that are used which, in turn, will significantly impact the efficiency of both schemes.

*Lack of simple, available implementations.* There are no public implementations of stateful PIR schemes. Additionally, no previous scheme implements their stateful PIR scheme as part of their experimental analysis. This means that the computational overheads of running existing schemes are either extrapolated from stateless PIR implementations, or remain unavailable. Having simple, small, and available implementations is a significant advantage when it comes to assessing the efficiency and security guarantees that are provided, during security and scientific auditing processes.

### 2.3 Overview of FrodoPIR

Figure 1 gives a diagrammatic overview of the following steps.



**Fig. 1.** An overview of FrodoPIR. In (1), the server compresses their database **DB** (represented as a matrix) into **M**, via multiplication with the global matrix **A** that is derived randomly from a public seed  $\mu$ . The client downloads  $(\mu, \mathbf{M})$ , and in (2) they preprocess a query and store  $(\mathbf{b}, \mathbf{c})$ , note that  $\mathbf{b}$  is an LWE sample and is thus randomly distributed. In the online phase, in (3), the client creates their query by adding an indicator value  $x$  to the  $i^{\text{th}}$  vector entry of  $\tilde{\mathbf{b}}$ . In (4), the server multiplies the client query vector with their DB matrix and return the result,  $\tilde{\mathbf{c}}$ . Finally, in (5), the client subtracts  $\mathbf{c}$  from  $\tilde{\mathbf{c}}$  — rounding the result to remove any error terms — and learns the  $i^{\text{th}}$  row of DB. The full scheme is given in Section 4.

1. In the offline phase, the server interprets the database as a matrix, applies a compression function to said matrix and makes the results available as global public parameters. This compression function shrinks the size of the database by a factor of approximately  $m/\lambda$ , where  $\lambda$  is the security parameter and  $m$  is the number of database elements. Thus, the size of the parameters is no longer linear in the size of the database.
2. The client downloads the public parameters, and can compute  $c$  sets of pre-processed query parameters.
3. In the online phase, the client uses a single set of preprocessed query parameters to produce an “encrypted” query vector, which is sent to the server.
4. The server responds to the query by multiplying the vector with their database matrix.
5. The client returns the result by “decrypting” the response using their pre-processed query parameters.

The security of the system relies on the decisional LWE problem<sup>3</sup>: each client query is a noisy vector that appears uniformly random to the server. Security conservatively holds up to a global number of client queries that the server witnesses. When this is reached, the server simply reruns the compression function using newly sampled randomness, and the clients download and process the new parameters.<sup>4</sup>

While the ideas behind FrodoPIR are fundamentally similar to previous RLWE-based PIR schemes, the key differentiating factor is that it uses a secure, client-independent preprocessing phase. Moreover, the total client download is much smaller than schemes that involve streaming the entire server database. This trade-off results in a scheme that is significantly cheaper than all previous approaches, including those that achieve sublinear computation and communication complexities such as [30]. FrodoPIR is especially well-suited to operating on databases containing *many* small elements. Moreover, FrodoPIR consistently achieves low runtimes across various database shapes (see Section 6).

The main limitation of the FrodoPIR approach is that online client queries are linear in the size of the database, which can be much larger than previous schemes. Fortunately, we show that FrodoPIR is highly configurable and that we are able to reduce client query sizes (as well as server-side online computation) at the cost of increasing the client download size (see Section 5.4 for more details). Another limitation is that the server database transformation can result in storing a larger amount than the database itself. Roughly speaking, the server database matrix is  $3\times$  as large as the original database. Such database transformations are common in PIR: RLWE-based schemes usually store their database in a format that allows using number-theoretic transform operations easily; and store database elements as FHE ciphertexts which can lead to a  $2\times$  increase in database storage.

We provide a functionality, efficiency, and coarse-level financial comparison between FrodoPIR and previous stateless/stateful PIR schemes in Table 1. We illustrate how amortization of offline computation across *all* client leads to significant efficiency advantages compared with previous stateful PIR schemes in Section 6.

### 3 Preliminaries

#### 3.1 Notation

We denote vectors and matrices in lower- and upper-case bold-face, respectively. All vectors  $\mathbf{v}$  are assumed to be in column orientation, and we write  $\mathbf{v}^T$  to denote the same vector in row orientation. For a set of vectors  $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ , we

<sup>3</sup> We base security specifically on a well-known variant of LWE, known as the decisional *ternary* LWE problem (Assumption 1).

<sup>4</sup> Less conservative security analyses suggest that the number of queries may not have such a strong impact on security, see Appendix C for wider discussion.



write  $[\mathbf{x}_1 \parallel \cdots \parallel \mathbf{x}_\ell]$  to denote the matrix with the  $i^{\text{th}}$  column set to equal  $\mathbf{x}_i$  for  $i \in [\ell]$ .

Let  $\lfloor x \rfloor \in \mathbb{Z}$  denote rounding  $x \in \mathbb{R}$  to the nearest integer, rounding down in case of a tie. Likewise, we use  $\lceil x \rceil$  to indicate explicitly rounding  $x \in \mathbb{R}$  to the next highest integer. For  $\mathbf{x} \in \mathbb{Z}_q^m$ , we write  $\lceil \mathbf{x} \rceil_p$  to denote the computation of  $\lceil p/q \cdot \mathbf{x} \rceil$ , where the rounding is applied to each entry of  $\mathbf{x}$  individually. For some set  $\mathcal{X}$ , we write  $x \leftarrow \mathcal{X}$  to denote that  $x$  is sampled from  $\mathcal{X}$  uniformly, and we write  $\mathbf{x} \leftarrow \mathcal{X}^m$  to denote sampling an  $m$ -dimensional vector, with each entry sampled uniformly from  $\mathcal{X}$ . We write  $\log(x)$  to denote taking the base-2 logarithm of  $x$ . We use  $\lambda$  to denote the security parameter throughout, and say that an algorithm  $\mathcal{A}$  is PPT if it runs in probabilistic polynomial-time with respect to  $\lambda$ .

### 3.2 Mathematical Preliminaries

We use the learning with errors (LWE) problem in its decisional version, which is equivalent to its search version as proven by Regev [66] for prime moduli, and was later shown to be equivalent for all moduli [60, 20].

**Definition 1.** (Decisional LWE problem) *Let  $\chi$  be some distribution, and let  $q, n, m > 0$  depend on  $\lambda$ . The decision LWE problem ( $\text{LWE}_{q,n,m,\chi}$ ) is to distinguish between:*

$$(\mathbf{A}, \mathbf{s}^T \cdot \mathbf{A} + \mathbf{e}^T) \in \mathbb{Z}_q^{n \times m} \times \mathbb{Z}_q^m, \quad (1)$$

$$(\mathbf{A}, \mathbf{u}) \in \mathbb{Z}_q^{n \times m} \times \mathbb{Z}_q^m, \quad (2)$$

where  $\mathbf{A} \leftarrow \mathbb{Z}_q^{n \times m}$ ,  $\mathbf{s}^T \leftarrow (\chi)^n$ ,  $\mathbf{e}^T \leftarrow (\chi)^m$ , and  $\mathbf{u} \leftarrow \mathbb{Z}_q^m$ .

Evidence that  $\text{LWE}_{q,n,m,\chi}$  is hard to solve for appropriate choices of  $\chi$  — for example small Gaussian distributions — and for both classical and quantum adversaries follows via reduction from standard worst-case lattice problems [66] (as hard as worst case problems on  $n$ -dimensional lattices).

*Variants of LWE.* The following assumption states that decisional  $\text{LWE}_{q,n,m,\chi}$  remains hard when  $\chi$  is the uniform distribution over ternary values (i.e.  $\{0, \pm 1\}$ ).

**Assumption 1** (Ternary LWE [20]) *The  $\text{LWE}_{q,n,m,\chi}$  problem is hard to solve when  $\chi$  is the uniform distribution on  $\{-1, 0, 1\}$  (i.e. the uniform ternary distribution).*

It follows from the work of Brakerski et al. [20] that Decisional LWE with ternary secrets is as hard as the problem investigated by Regev [66]. Moreover, many examples of well-established cryptographic schemes rely on the hardness of LWE using both secrets and errors sampled from a uniform ternary distribution, such as [51, 15, 34, 46].

In Definition 2 we give a variant of  $\text{LWE}_{q,n,m,\chi}$  known as the Matrix LWE problem (denoted by  $\text{MatLWE}_{q,n,m,\chi,\ell}$ ). Corollary 1 shows that  $\text{MatLWE}_{q,n,m,\chi,\ell}$  is hard to solve, with only polynomial security loss compared with  $\text{LWE}_{q,n,m,\chi}$  [17].

**Definition 2.** (Decisional Matrix LWE problem [17]) Let  $\chi$  be some distribution, and let  $q, n, m, \ell > 0$  depend on  $\lambda$ . The decisional Matrix LWE problem ( $\text{MatLWE}_{q,n,m,\chi,\ell}$ ) is to distinguish between:

$$(\mathbf{A}, \mathbf{S} \cdot \mathbf{A} + \mathbf{E}) \in \mathbb{Z}_q^{n \times m} \times \mathbb{Z}_q^{\ell \times m}, \quad (3)$$

$$(\mathbf{A}, \mathbf{U}) \in \mathbb{Z}_q^{n \times m} \times \mathbb{Z}_q^{\ell \times m}, \quad (4)$$

where  $\mathbf{A} \leftarrow_{\$} \mathbb{Z}_q^{n \times m}$ ,  $\mathbf{S} \leftarrow (\chi)^{\ell \times n}$ ,  $\mathbf{E} \leftarrow (\chi)^{\ell \times m}$ , and  $\mathbf{U} \leftarrow_{\$} \mathbb{Z}_q^{\ell \times m}$ .

**Corollary 1.** (Hardness of  $\text{MatLWE}_{q,n,m,\chi,\ell}$  [17]) Let  $\mathcal{A}$  be a PPT adversary against  $\text{MatLWE}_{q,n,m,\chi,\ell}$  with advantage  $\epsilon$ , then we can construct a PPT adversary  $\mathcal{B}$  against  $\text{LWE}_{q,n,m,\chi}$  with advantage  $\epsilon/\ell$ .

We now state the following as a corollary of the central limit theorem, to provide an upper bound on the size of sums of elements sampled from uniform ternary distributions.

**Corollary 2.** (Bounds on uniform ternary sums) For sufficiently large  $m = \text{poly}(\lambda)$ , the sum of  $m$  samples taken from the uniform distribution over ternary values (i.e.  $\{-1, 0, 1\}$ ) is bounded by  $4\sqrt{m}$  with all but negligible probability.

### 3.3 Stateful Private Information Retrieval

As discussed, in this work, we will consider *stateful* PIR schemes, where the PIR interactions are split into a query-independent offline phase and a query-dependent online phase [65]. A stateful PIR scheme consists of an offline and an online phase, which are defined as follows.

**Offline phase:**

- $\text{ssetup}(1^\lambda)$ : An algorithm run by the server that outputs some initialization parameters  $\text{ip}$ .
- $\text{cinit}(\text{ip})$ : A client initialization algorithm run on parameters  $\text{ip}$ . Outputs a message  $\text{msg}$  to be sent to the server during the offline phase.
- $\text{spreproc}(\text{ip}, \text{DB}, \text{msg})$ : A server preprocessing algorithm run on  $\text{ip}$ , the server database  $\text{DB}$ , and client message  $\text{msg}$ . Outputs a set of public parameters  $\text{pp}$  to be downloaded by the client.
- $\text{cpreproc}(\text{ip}, \text{pp})$ : A client preprocessing algorithm run on the server-generated public parameters  $(\text{ip}, \text{pp})$ , that outputs a state  $\text{st}$ .

Stateful PIR schemes that omit the  $\text{cinit}$  algorithm are said to have *client-independent* preprocessing phases.

**Online phase:**

- $\text{query}(\text{st}, i)$ : A client algorithm that generates a PIR query  $\mathbf{q}$  for the item in the  $i^{\text{th}}$  index of the server database, and optionally returns an updated state  $\text{st}'$ .
- $\text{respond}(\text{DB}, \mathbf{q})$ : A server algorithm that outputs a response  $\mathbf{r}$  to be returned to the client.
- $\text{process}(\text{st}, \mathbf{r})$ : A client algorithm that takes the sever response  $\mathbf{r}$ , and outputs a database element  $\mathbf{x}$ .

### 3.4 PIR requirements

Stateful PIR schemes must guarantee the following.

*Correctness.* The following correctness definition ensures that clients receive the correct response with overwhelming probability when interacting with an honest server.

**Definition 3.** (Correctness) *Let  $\text{DB}$  be the server database, let  $i$  be the index that the client seeks to query during the online phase, and let  $\text{DB}[i]$  be the  $i^{\text{th}}$  entry of  $\text{DB}$ . We say a PIR scheme is correct if the following inequality is satisfied.*

$$\Pr \left[ x = \text{DB}[i] \begin{array}{l} \text{ip} \leftarrow \text{ssetup}(1^\lambda) \\ \text{pp} \leftarrow \text{spreproc}(\text{ip}, \text{DB}, \text{cinit}(\text{ip})) \\ \text{st} \leftarrow \text{cpreproc}(\text{ip}, \text{pp}) \\ \text{q} \leftarrow \text{query}(\text{st}, i) \\ \text{r} \leftarrow \text{respond}(\text{DB}, \text{q}) \\ \text{x} \leftarrow \text{process}(\text{st}, \text{r}) \end{array} \right] > 1 - \text{negl}(\lambda)$$

*Security.* We use the standard definition of security in enforcing the indistinguishability of client queries. As is common throughout PIR literature, this assumes a semi-honest server, that follows the protocol correctly and attempts to learn more based on the client messages they receive.

**Definition 4.** (1-query indistinguishability) *As stated, let  $\text{DB}$  be the server database. Generate the server public parameters by running  $\text{ip} \leftarrow \text{ssetup}(1^\lambda)$  and  $\text{pp} \leftarrow \text{spreproc}(\text{ip}, \text{DB}, \text{cinit}(\text{ip}))$ , and let  $\text{st} \leftarrow \text{cpreproc}(\text{ip}, \text{pp})$  be the client state. We say that a PIR scheme is secure if, for any PPT adversary  $\mathcal{A}$  specifying indices  $i_0, i_1$  that is given  $\text{q}_b \leftarrow \text{query}(\text{st}, i_b)$  for  $b \leftarrow_{\$} \{0, 1\}$ , then  $\mathcal{A}$  has negligible probability in correctly guessing  $b$ .*

The above definition can be expanded to specify  $\ell$ -query indistinguishability, in other words that two sets of size  $\ell$  of client queries are indistinguishable from each other [65].

*Efficiency.* PIR schemes require a communication overhead smaller than the solution of having clients download the entire server database. In the stateful PIR case, it should hold when amortizing costs over the number of client queries.

**Definition 5.** (Efficiency) *For a single client launching  $c$  queries, a PIR scheme is efficient if the total client communication overhead is smaller than  $|\text{DB}|$ .*

Therefore, for stateful schemes, the total client communication cost is calculated using the equation:  $\text{comms}(\text{offline}) + c \cdot \text{comms}(\text{online})$ .

## 4 Our Scheme

We now describe the FrodoPIR scheme, writing FPIR for short.

## 4.1 Cryptographic Setup

Recall that  $\mathcal{S}$  is the server holding the database DB that each client attempts to learn entries from. DB is an array of  $m$  elements, each made up of  $w$  bits. Each entry is associated with the index  $i$  that corresponds to its position in the array. For now, we will assume that the client knows which index it would like to query during the online phase of the protocol.<sup>5</sup> We assume that there are  $C$  clients that will each launch a maximum of  $c$  queries against DB. Regarding the LWE instantiation that is used: let  $n$  and  $q$  be the LWE dimension and modulus, respectively; let  $\rho$  be the number of bits packed into each entry of the DB matrix; let  $0 < \rho < q$ ; let  $\chi$  be the uniform distribution over  $\{-1, 0, 1\}$ ; and let  $\lambda$  be the concrete security parameter. Finally, we use  $\text{PRG}(\mu, x, y, q)$  to denote a pseudorandom generator that expands a seed  $\mu \leftarrow_{\$} \{0, 1\}^\lambda$  into a matrix in  $\mathbb{Z}_q^{x \times y}$ .

## 4.2 Preprocessing Phase

We first describe the offline phase which occurs before the client makes any queries to the server. Note that `cinit` is not required in FrodoPIR, and thus we do not define it.

**Server setup (FPIR.ssetup).** The server constructs their database containing  $m$  elements, and samples a seed  $\mu \in \{0, 1\}^\lambda$ .

**Server preprocessing (FPIR.spreproc).** In this phase, the server derives a matrix  $\mathbf{A} \leftarrow \text{PRG}(\mu, n, m, q)$ , where  $\mathbf{A} \in \mathbb{Z}_q^{n \times m}$ . It then runs  $\mathbf{D} \leftarrow \text{parse}(\text{DB}, \rho)$ , where `parse` encodes the DB into a matrix  $\mathbf{D} \in \mathbb{Z}_\rho^{m \times \omega}$ , and where  $\omega = \lceil w / \log(\rho) \rceil$ .<sup>6</sup> The server stores  $\mathbf{D}$ .

To generate public parameters, the server runs  $\mathbf{M} \leftarrow \mathbf{A} \cdot \mathbf{D}$ , and then publishes the pair  $(\mu, \mathbf{M}) \in \{0, 1\}^\lambda \times \mathbb{Z}_q^{n \times \omega}$  to a public repository accessible by clients.

**Client preprocessing (FPIR.cpreproc).** Each client downloads  $(\mu, \mathbf{M})$  from the public repository, and runs  $\mathbf{A} \leftarrow \text{PRG}(\mu, n, m, q)$ . The client then samples  $c$  vectors  $\mathbf{s}_j \leftarrow (\chi)^n$  and  $\mathbf{e}_j \leftarrow (\chi)^m$ . Finally, it computes  $\mathbf{b}_j \leftarrow \mathbf{s}_j^T \cdot \mathbf{A} + \mathbf{e}_j^T \in \mathbb{Z}_q^m$  and  $\mathbf{c}_j \leftarrow \mathbf{s}_j^T \cdot \mathbf{M} \in \mathbb{Z}_q^\omega$ , for each  $j \in [c]$ , and stores the pairs internally as the set  $X = (\mathbf{b}_j, \mathbf{c}_j)_{j \in [c]}$ .

## 4.3 Online Phase

**Client query generation (FPIR.query).** For the index  $i$  that the client wishes to query, the client generates a vector  $\mathbf{f}_i = (0, \dots, 0, q/\rho, 0, \dots, 0)$ , i.e. the all-zero vector except where  $\mathbf{f}_i[i] = q/\rho$ . The client then pops a pair  $(\mathbf{b}, \mathbf{c})$  from the internal state `st` that it maintains, and computes  $\tilde{\mathbf{b}} \leftarrow \mathbf{b} + \mathbf{f}_i \in \mathbb{Z}_q^m$ , and sends  $\tilde{\mathbf{b}}$  to the server.

<sup>5</sup> Section 7 discusses options for mapping string-based queries to indices.

<sup>6</sup> Thus, the  $i^{\text{th}}$  row consists of  $\omega \log(\rho)$ -bit chunks of  $\text{DB}[i] \in \mathbb{Z}_\rho^\omega$ .

**Server response (FPIR.respond).** The server receives  $\tilde{\mathbf{b}}$  from the client, and responds with  $\tilde{\mathbf{c}} \leftarrow \tilde{\mathbf{b}}^T \cdot \mathbf{D} \in \mathbb{Z}_q^\omega$ .

**Client postprocessing (FPIR.process).** The client receives  $\tilde{\mathbf{c}}$ , and outputs  $\mathbf{x} \leftarrow \lfloor \tilde{\mathbf{c}} - \mathbf{c} \rfloor_\rho \in \mathbb{Z}_\rho^\omega$ .

#### 4.4 Correctness

The output of the client postprocessing phase is  $x \leftarrow \lfloor \tilde{\mathbf{c}} - \mathbf{c} \rfloor_\rho$ . Expanding the right-hand side of the equation gives:

$$\begin{aligned} \mathbf{x} &= \lfloor \tilde{\mathbf{c}} - \mathbf{c} \rfloor_\rho \\ &= \lfloor (\mathbf{s}^T \cdot \mathbf{A} + \mathbf{e}^T + \mathbf{f}_i^T) \cdot \mathbf{D} - (\mathbf{s}^T \cdot \mathbf{A} \cdot \mathbf{D}) \rfloor_\rho \\ &= \lfloor (\mathbf{e} + \mathbf{f}_i)^T \cdot \mathbf{D} \rfloor_\rho. \end{aligned} \quad (5)$$

Noting that  $\lfloor \mathbf{f}_i^T \cdot \mathbf{D} \rfloor_\rho = \mathbf{D}[i]$  where the  $i^{\text{th}}$  row  $\mathbf{D}[i] \in \mathbb{Z}_\rho^\omega$  is interpreted as a column vector, then proving that

$$\lfloor \mathbf{e}^T \cdot \mathbf{D} + \mathbf{f}_i^T \cdot \mathbf{D} \rfloor_\rho = \lfloor \mathbf{f}_i^T \cdot \mathbf{D} \rfloor_\rho \quad (6)$$

results in the client learning the correct output. This gives rise to the following theorem.

**Theorem 2.** (Correctness) *If  $q \geq 8\rho^2\sqrt{m}$ , then FPIR is correct with high probability.*

*Proof.* See Appendix A.1.

#### 4.5 Security

To prove security of FrodoPIR, we show that any query  $\tilde{\mathbf{b}} \leftarrow \text{FPIR.query}(i)$  is distributed uniformly in  $\mathbb{Z}_q^m$  from the perspective of  $\mathcal{S}$  (Theorem 3). This general result proves that FPIR satisfies 1-query indistinguishability (Definition 4) and we further show that this holds for  $\ell = \text{poly}(\lambda)$  client queries in Corollary 3. Since,  $\chi$  is the uniform ternary distribution, the proofs therefore follow from the hardness of the decisional ternary LWE problem (Assumption 1).

**Theorem 3.** (1-query indistinguishability) *FPIR is secure under observation of 1 query, under the assumption that  $\text{LWE}_{q,n,m,\chi}$  is difficult to solve.*

*Proof.* See Appendix A.2.

**Corollary 3.** ( $\ell$ -query indistinguishability) *FPIR is secure under observation of  $\ell = \text{poly}(\lambda)$  queries, under the assumption that  $\text{MatLWE}_{q,n,m,\chi,\ell}$  is difficult to solve.*

*Proof.* See Appendix A.3.

**Table 2.** Communication overheads (bits) of FrodoPIR.

	Offline	Online
Client upload	—	$m \log(q)$
Client download	$\lambda + n\omega \log(q)$	$\omega \log(q)$

**Table 3.** Number of operations required in FrodoPIR.

	spreproc	cpreproc	query	respond	process
mod $q$ mults	$nm\omega$	$n(m + \omega)$	—	$m\omega$	—
mod $q$ adds	$n(m - 1)\omega$	$(n - 1)(m + \omega)$	1	$(m - 1)\omega$	$\omega$
PRG	$nm$	$nm$	—	—	—

#### 4.6 Efficiency

We give the conditions under which FPIR satisfies the efficiency goal of a PIR scheme, as laid out in Definition 5.

**Theorem 4.** (Efficiency) *Let  $c$  be the upper bound of a single client’s FPIR queries. Then FPIR is efficient when:*

$$\lambda + n\omega \log(q) + c(m + \omega) \log(q) < |DB|.$$

*Proof.* This follows from Definition 5, considering the communication costs of FrodoPIR (see Table 2).

## 5 Parameter Settings and Configurations

We now describe parameter settings and potential optimizations that demonstrate the versatility of FrodoPIR. The major parameters of the scheme to be configured are: the concrete security parameter  $\lambda$ ; the LWE dimension  $n$ ; the LWE modulus  $q$ ; the uniform ternary distribution,  $\chi$ , used for sampling LWE secret and error vectors; the number of bits,  $\rho$ , packed into each entry of the DB matrix,  $\mathbf{D}$ ; the number of elements,  $m$ , in the server DB; and the bit-length,  $w$ , of each element in the server database.

The communication overheads of FrodoPIR are given in Table 2, the number of required computational operations are given in Table 3, and the storage overheads in Table 4.<sup>7</sup> Clearly, the aforementioned parameters all have an impact on the protocol efficiency.

<sup>7</sup> Recall that  $\omega = w/\log(\rho)$ .

**Table 4.** Storage overheads of FrodoPIR in bits, according to whether client performs any offline preprocessing of queries (where  $c$  is the number of preprocessed queries), or not. When no preprocessing is performed, the client storage overhead is logarithmically dependent on the number of elements in DB.

	with preprocessing	without
Server storage	$\lambda + m\omega \log(\rho)$	$\lambda + m\omega \log(\rho)$
Client storage	$\lambda + c(m + \omega) \log(q)$	$\lambda + n\omega \log(q)$

### 5.1 Required Invariants

Firstly, for efficiency, FrodoPIR must satisfy Theorem 4:

$$\lambda + n\omega \log(q) + c(m + \omega) \log(q) < mw. \quad (7)$$

Secondly, for correctness (Theorem 2), we must have that:

$$q \geq 8\rho^2 \sqrt{m}, \quad (8)$$

holds. Finally, for security,  $\text{MatLWE}_{q,n,m,\chi,\ell}$  must provide at least 128 bits of concrete classical security. We can estimate the concrete security of LWE instances with the lattice security estimator [2], see Appendix D.

### 5.2 Fixing LWE Parameters

Before configuring FrodoPIR for efficiency, we first fix a set of parameters that provide the necessary concrete security guarantees. We focus on those parameters for  $\text{MatLWE}_{q,n,m,\chi,\ell}$ , except for  $m$  which is the number of DB elements.

Firstly,  $\chi$  is the uniform ternary distribution. Secondly, we set  $q = 2^{32}$ , which allows us to use standard 32-bit unsigned integer operations for the implementation of the  $\mathbb{Z}_q$  operations. Thirdly, we set  $n = 1774$  as the LWE dimension. This choice conservatively estimates the security of the  $\text{MatLWE}_{q,n,m,\chi,\ell}$  instance, using the work of Albrecht et al. [2] to provide the security of  $\nu = \text{LWE}_{q,n,m,\chi}$ , and then calculate the concrete security parameter as  $\lambda = \nu - \log \ell$  (Corollary 1) using the primal-USVP cost model. As  $\ell$  is the total number of queries that the server observes, we set  $\ell = 2^{52}$  queries as a suitable upper bound before rotation of  $\mathbf{A}$  is required. When  $\mathbf{A}$  is rotated, the security of the scheme is reset. Therefore,  $\lambda = \nu - 52$ . The code that we eventually run for estimating the security of  $\nu$  is given in Appendix D.

The conservative analysis above dictates that for larger numbers of queries, the concrete security of the instance will decrease polynomially in the number of queries — until a new LWE matrix  $\mathbf{A}$  is resampled.<sup>8</sup> Note that no attacks currently exist that exploit the security of  $\text{MatLWE}_{q,n,m,\chi,\ell}$  in this way. As such,

<sup>8</sup> Since PIR is constructed in a semi-honest security model, we safely assume that the server resamples  $\mathbf{A}$  when it is required to do so.

**Table 5.** Database, query, and security parameter settings.

$q$	$2^{32}$	$2^{32}$	$2^{32}$	$2^{32}$	$2^{32}$
$n$	1774	1774	1774	1774	1774
$m$	$2^{16}$	$2^{17}$	$2^{18}$	$2^{19}$	$2^{20}$
$\rho$	$2^{10}$	$2^{10}$	$2^{10}$	$2^9$	$2^9$
$\kappa$	13.028	26.056	52.112	93.802	187.603
$\lambda$	128	128	128	128	128

a less conservative analysis may be valuable in allowing smaller dimensions  $n$ , by simply estimating the security for smaller values of  $\ell$ , or by estimating the security of  $\nu = \text{LWE}_{q,n,m,\chi}$  only. We discuss the performance impact of choosing smaller values of  $n$  in Appendix C.

### 5.3 Recommended Database Parameters

Let  $\kappa = (\log(\rho) \cdot m)/(n \cdot \log(q))$  denote the improvement factor relative to the offline client download when compared to the original DB size. In Table 5, we give recommended parameter settings for FrodoPIR. For each parameter set, the concrete security parameter is 128 bits for  $2^{52}$  client queries. Security can be increased by increasing the dimension  $n$ , though, this reduces  $\kappa$ . The lattice security estimates that we derive are produced using the code detailed in Appendix D.

In Section 6, we consider DB elements of size  $w = 1\text{KB}$ , which leads to  $\omega \in \{820, 911\}$ , depending on the value of  $\rho$ . In Section 6.3, we also highlight how performance changes when considering DB elements of larger sizes.

### 5.4 Additional Optimizations

*Processing larger databases via sharding.* As  $m$  increases beyond  $2^{20}$ , we see a greater relative saving of download costs relative to the fixed  $n$  that is used (as parameterized by  $\kappa$ ). However, this has undesirable impacts on the performance of the scheme. First, all online server-side computation in the online phase is linearly dependent on  $m$ , and so increasing  $m$  immediately results in higher latency. The offline work scales similarly for client devices, which are typically constrained and unlikely to cope with vast overheads. Second, the client query size rapidly grows as it is also linearly dependent on  $m$ .

Overall, we expect that the best approach for operating on larger databases is to *shard* them into  $s$  parallel instances, each using a database of size  $m/s$ . Each instance can then independently process the *same* single client query. This allows the client to learn the  $i^{\text{th}}$  index from each of the  $s$  shards from only a single query. This allows parallelization of server computation, and careful management of available computing resources. On the client-side, the size of the online query is linear in  $m/s$ , rather than  $m$ , which can lead to more efficient uploads. However, this comes at the cost of increasing the client download by a



factor of  $s$ . As a result, sharding allows developing different trade-offs for client upload/download for various situations. Previous work has already highlighted the benefits of performing such sharding on the server database [33] in terms of increasing amortization factors and allowing further degrees of parallelization.

Note that each client must download the public parameters of each of the individual shards. This increases the size of the client download, but with the benefits of reducing the size of their own query and reducing server-side latency. Additionally, noting the independence of each server-side vector-column multiplication in FrodoPIR, we could equally leverage sharding by splitting the server database matrix into smaller subsets of columns for handling larger data elements.

*Database updates.* Sharding alone does not reduce the client overhead in preprocessing queries, which remain linear in the total database size. This can become expensive if the server database is updated frequently: each time the client has to regenerate their preprocessed query data.

However, coupling sharding with a database updating procedure that touches only few of the shards can reduce database updates to only re-running the `ssetup`, `cpreproc`, and `spreproc` procedures on a small fraction of the database. Specifically, if database updates are confined to a single shard of the database, then these procedures need only be run on that particular shard after every update. Updating a single shard of the database results in only requiring the client to download and process an amount of data that is a  $1/(\kappa \cdot s)$  fraction of the entire database. Even for large databases, this fraction is likely to be very small.

*Achieving logarithmic client-storage overhead.* Table 4 clearly highlights that storage overheads for clients are dependent on  $c$ , the number of preprocessed queries. These costs can be reduced significantly to being logarithmically dependent on  $m$ , by simply not performing any preprocessing. The reason that the costs are logarithmic is that the client storage is equal to  $(\lambda + n\omega \log(q))$  where, as mentioned in Section 5,  $q$  is chosen to be equal to  $8\rho^2\sqrt{m}$ . This approach requires derivation of the matrix  $\mathbf{A}$  and query parameters for every online query. Since the derivation of  $\mathbf{A}$  is fairly costly, computation-constrained clients will benefit from preprocessing client queries.

## 6 Experimental Analysis

We provide an experimental analysis of the incurred computational runtimes, bandwidth usage, and financial costs when running FrodoPIR. Further, we highlight how such costs amortize over the one-time offline preprocessing phase. Finally, we compare these costs with the previous stateful PIR schemes — PSIR [65], SOnionPIR [62], and CHKPIR [30].

*Performance benchmarks.* We run all experiments as single-threaded processes on an Amazon `t2.2xlarge` EC2 instance, with 8 CPU cores and 32GB of RAM.<sup>9</sup>

<sup>9</sup> Client-based functions are estimated using the same hardware.

This is equivalent to the setup that is used in [62] for comparing SOnionPIR and PSIR. All computational costs correspond to CPU processing time. Bandwidth costs are calculated using the overheads detailed in Table 2, where  $\lambda = 128$ . Regarding financial cost calculations, transferring data from server to clients costs \$0.09 per GB, the cost of data transfer from clients to server is free, and the cost of computation is \$0.3712 per hour of usage (or \$0.0464 per CPU hour).<sup>10</sup>

*Parameter choices.* We provide non-amortized communication and computation benchmarks for a single server database using each of the parameter settings provided in Table 5. We choose  $w = 2^{13}$  bits (i.e. 1KB database elements); and set  $\rho = 2^{10}$  for  $m \leq 2^{18}$ , and  $\rho = 2^9$  otherwise. Section 6.3 provides additional benchmarks for different database shapes, such as in cases where database elements are much larger.

The parameters we use conservatively provide 128-bit security for around  $2^{52}$  client queries. In Appendix C, we discuss performance improvements when the lattice dimension,  $n$ , is reduced to account for less conservative security estimations.

*Source code.* Our open-source<sup>11</sup> implementation of FrodoPIR is written in Rust, consists of 735 lines of code, including tests, and requires no external dependencies relating to cryptographic operations. All modular arithmetic is implemented using instructions associated with the 32-bit unsigned integer type included in the Rust standard library.

*Example application.* In Appendix B, we further illustrate how FrodoPIR can be applied to real-world use-cases, taking, as an example, the Google SafeBrowsing API [44].

## 6.1 Performance Results

Table 6 displays the non-amortized performance of FrodoPIR. This involves calculating running times and bandwidth usage for running a single-threaded server instance in interaction with a single client. Later, we analyze how the offline costs amortize on a per-query basis. Amortization is calculated over a variable number of clients  $C$ , and the number of per-client queries  $c$  (where we set  $c = 500$  for all experiments).

*Offline phase.* The server generates their database matrix DB and public parameters. This is a client-independent operation that scales linearly in  $m$ . This process includes pseudorandom derivation of the LWE matrix  $\mathbf{A} \in \mathbb{Z}_q^{n \times m}$  from a single  $\lambda$ -bit seed, which is also computed by each client. After downloading the public parameters, the client performs query-independent preprocessing for each query that they will make. The results of preprocessing are used during the online phase. These costs grow roughly linearly in  $m$ .

<sup>10</sup> <https://aws.amazon.com/ec2/pricing/on-demand/>, August 2022.

<sup>11</sup> <https://github.com/brave-experiments/frodo-pir>

**Table 6.** Non-amortized performance analysis of FrodoPIR. The “Client derive matrix” cost refers to the cost of deriving the LWE matrix  $\mathbf{A}$  from the seed  $\mu$ , while “Client query preprocessing” refers to the cost of query-independent preprocessing required for a single query. The server offline phase costs can be amortized *globally* across the number of queries ( $C$ ) that are performed, while the client download and parameter derivation costs amortizes across the number of queries ( $c$ ) that they individually make.

	Number of DB items ( $\log(m)$ )	16	17	18	19	20
Offline	Client download (KB)	5682.47	5682.47	5682.47	6313.07	6313.07
	Database preprocessing (s)	104.57	206.26	429.07	936.36	1895.2
	Client derive matrix (s)	0.5826	1.1698	2.2118	4.7284	9.25
	Client query preprocessing (s)	0.1468	0.2898	0.5795	1.182	2.343
Online	Client query (KB)	256	512	1024	2048	4096
	Server response (KB)	3.203	3.203	3.203	3.556	3.556
	Client query (ms)	0.0213	0.0422	0.0811	0.1648	0.3429
	Server response (ms)	45.013	94.505	188.36	417.92	825.37
	Client output (ms)	0.359	0.398	0.363	0.42	0.434

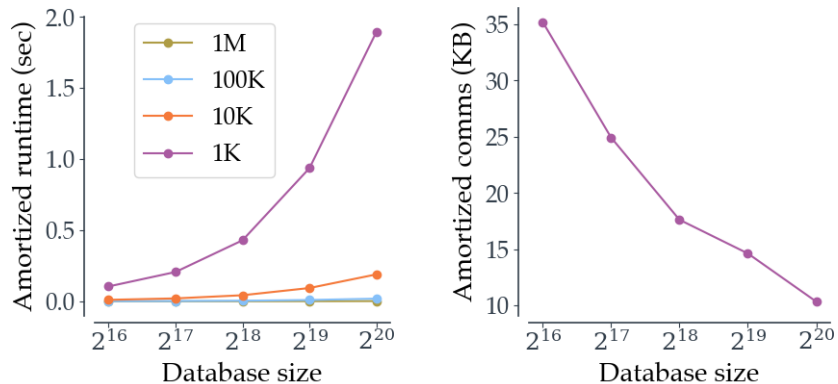
In terms of communication, the server publishes the  $\lambda$ -bit seed,  $\mu$ , and the matrix  $\mathbf{M} \in \mathbb{Z}_q^{n \times \omega}$ , where  $\omega = w/\log(\rho)$ . The size of the client download is static for each choice of  $\log(\rho)$ . As a consequence, the total cost only grows when increasing  $m$  dictates that  $\rho$  must also decrease.

*Online phase.* The client computation consists of performing a single addition operation to modify a single portion of preprocessed data. The client also performs a very small amount of postprocessing of servers responses that is almost static across all experiments, as it is linear in  $\omega$ . The dominant computation cost is the server-side processing of the client query that is  $\leq 0.83$ s for all database sizes.

The dominant communication cost relates to the client query, which is equal to  $m \log(q)$  bits and scales linearly in the DB size. The server response is significantly smaller —  $\omega \log(q)$  bits — resulting in a  $< 3.6\times$  overhead in the server response size compared with the original 1KB data element.

*Amortization of offline phase.* Many of the offline costs in Table 6 can be amortized significantly over the number of queries that are launched. In Figure 2, we give an overview of the rate of this amortization for DB preprocessing and parameter generation steps (when servicing between 1K and 1M queries), as well as the cost of the client downloads. The expensive cost of the one-time preprocessing of DB amortizes over all queries *globally*, i.e. over all clients. The client offline preprocessing and download amortizes over the value of  $c$ .

The total amortized computation cost (per-query) for the server and clients are given in Figure 3. We display server offline costs that are amortized across all client queries globally for between 1K and 1M queries. The majority of server



**Fig. 2.** Amortized (per-query) cost of server preprocessing (**left**), according to how many client queries they service, and client offline download size (**right**).

costs occur during the relatively cheap online phase. The majority of client work is performed during the query-independent offline phase, part of which (the derivation of  $\mathbf{A}$ ) can be amortized over  $c$ . Online costs for clients are very small.

*Storage costs.* Figure 4 illustrates the growth of the client storage overhead associated with the database size, when preprocessing  $c$  queries during the offline phase. This becomes fairly large when  $|\text{DB}| = 2^{20}$ , equalling roughly 2GB.

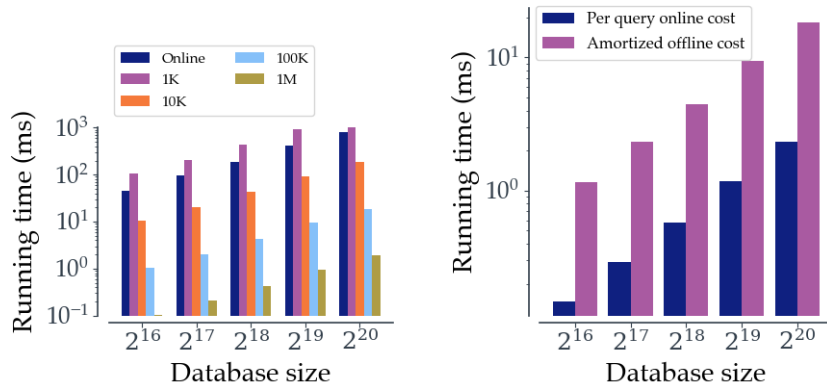
As mentioned in Section 5.3, it is possible to achieve  $\log(m)$  client-side storage overhead, which may be valuable for storage-constrained clients. The downside is that client online query processing grows noticeably, as seen in the right-side of Figure 4. This is due to having to perform all query-related processing in the online phase, including the derivation of  $\mathbf{A}$  from the public parameters (which can take from between 0.5 to 9.25 seconds, depending on the database size).<sup>12</sup>

*Financial costs.* The server-side financial costs given in Figure 5 take into account the expenses associated with both bandwidth and single-threaded computation. The initial preprocessing of the server database is a little higher than 1 cent for a database of  $2^{20}$ . The online per-query cost is tiny in comparison, and approximately 0.001 of a cent even for the largest DB size. The total per-query cost is calculated as the amortized offline costs, plus the online per-query cost.

## 6.2 Comparison with Prior Work

In Figure 6, we compare the performance of FrodoPIR with SOnionPIR [62] and PSIR [65]. Our comparison focuses on three performance criteria: computational

<sup>12</sup> The matrix  $\mathbf{A}$  must be rederived on usage to achieve  $\log(m)$  storage.



**Fig. 3.** Total online and amortized (per-query) offline computation costs for the server (**left**), according to how many client queries they service, and for the client (**right**).

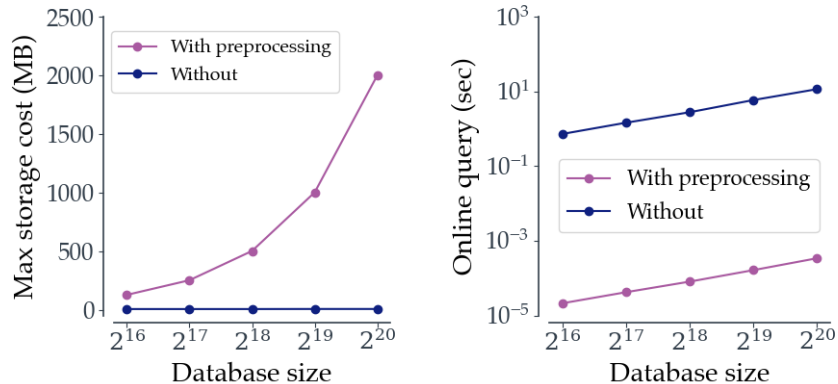
runtimes, bandwidth usage, and financial cost. Each comparison includes the cost of answering queries in FrodoPIR against the estimated<sup>13</sup> costs of running both SOnionPIR and PSIR.<sup>14</sup> Note that the costs presented in [62] result from estimating SOnionPIR and PSIR on the same EC2 hardware (`t2.2xlarge`) that we used for implementing FrodoPIR. We also provide details on how these costs amortize as the number of clients grows.

Our comparison considers total database sizes of  $|\text{DB}| \in \{2^{16}, 2^{18}, 2^{20}\}$ , and element sizes of 1KB. Note that SOnionPIR and PSIR allow packing of 30KB and 3KB of data into each server response [62]. This effectively allows shrinking the server DB by a factor of  $30\times$  and  $3\times$ , respectively, in kind. Since such costs are linear in the size of DB, we reduce the previously estimated runtime costs of both schemes accordingly. Offline costs for SOnionPIR are dependent on the number of queries,  $c$ , that are made by each client. For each DB size we set  $c = 500$ , the same value as used in [62]. For the financial costs, we provide costs per CPU hour of server-side computation. The comparison does not cover storage costs or client computation as neither measurement is explicitly provided by the previous schemes.

*Supporting databases with more elements.* Note that [62] provides estimated costs of the SOnionPIR and PSIR schemes for a DB of size  $2^{24}$ , but RAM overheads of FrodoPIR mean that the `t2.2xlarge` EC2 instance is not powerful enough to process a database of this size. This is also likely to be the case for the previous schemes. Building an efficient implementation of FrodoPIR for a database of  $2^{24}$  is possible by sharding, using 16 DB instances with  $2^{20}$  elements. In the interest of maintaining a fair comparison without using parallelization, we do not modify

<sup>13</sup> Neither previous stateful scheme has been fully implemented.

<sup>14</sup> Where PSIR uses SealPIR as the underlying PIR scheme.



**Fig. 4. Left:** Storage costs for clients demonstrating the trade-off between amortization of offline preprocessing and ensuring logarithmic storage overhead relative to  $m$ . **Right:** Comparison of online query costs when preprocessing, against performing all query-related computation in the online phase.

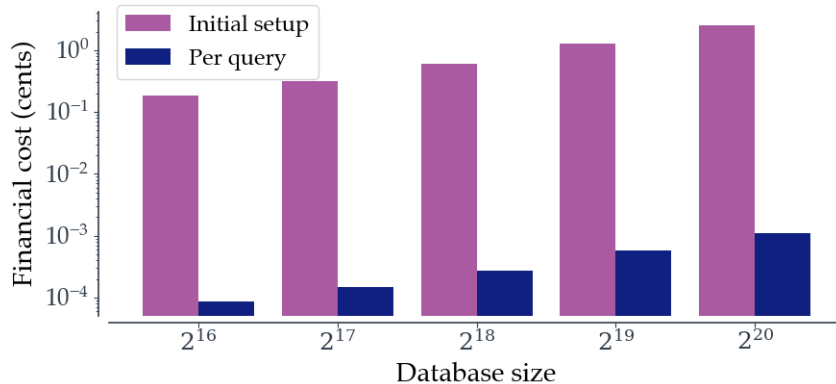
the hardware used or make use of sharding. Thus, we limit the comparison to database sizes  $\leq 2^{20}$ .

*Security levels.* We do not modify the security parameters of either SONionPIR or PSIR: they both offer  $\leq 115$  bits of security according to [2]. In contrast, FrodoPIR offers 128-bit security for up to 1 billion client queries and higher security levels for lower numbers. SONionPIR and PSIR could achieve higher security levels by doubling  $n$ ,<sup>15</sup> but while computation times would go unchanged, the server online response size would increase dramatically.

*Computation.* In the offline phase (Figure 6 (1)), the server-side computation for PSIR is zero, since the client simply downloads the entire server DB. The overall cost of computation in FrodoPIR grows linearly in the database size. While SONionPIR appears to outperform FrodoPIR for a single client, this cost increases linearly in the number of queries that a client wishes to make. As a consequence, if the number of queries per-client ( $c$ ) increases, then the cost of SONionPIR will quickly become greater. More importantly, as the number of clients ( $C$ ) in the system grows, this cost will continue to increase. In contrast, all preprocessing in FrodoPIR is client-independent, and thus it is fixed regardless of both  $c$  and  $C$ . Therefore, in a large multi-client deployment, it is clear that FrodoPIR is much cheaper than SONionPIR.

In the online phase (Figure 6 (2)), PSIR provides the fastest computation times. Both FrodoPIR and SONionPIR still provide competitive runtimes.

<sup>15</sup> Smaller  $n$  would suffice, but  $n$  has to be a power-of-two to ensure the efficiency of NTT-related optimizations.



**Fig. 5.** Financial costs (cents) associated with running the server in FrodoPIR. The initial setup cost can be amortized globally across all client queries.

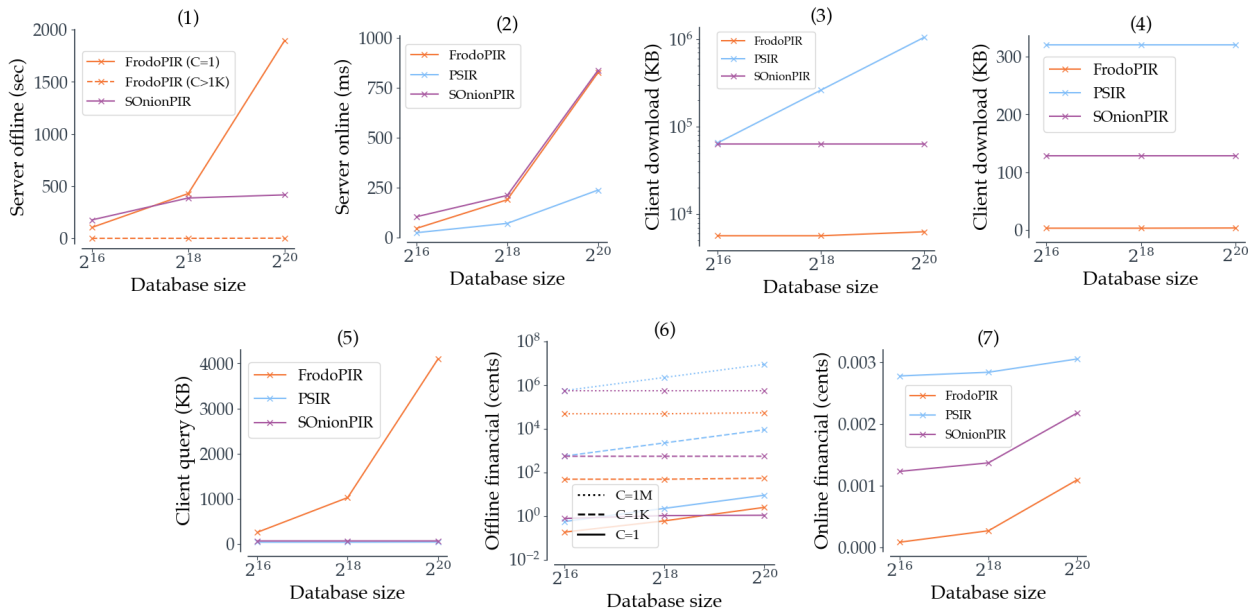
FrodoPIR requires  $\leq 0.825$ s for responding to a client query on a DB with  $2^{20}$  elements.

*Communication.* The offline client download cost (Figure 6 (3)) in SOnionPIR is heavily dependent on the number of queries that will be launched. The cost of PSIR is essentially the cost of downloading the entire server DB. Note that the client download in FrodoPIR grows logarithmically in the size of DB. Overall, since the costs of FrodoPIR amortize across the number of queries launched by the clients, with a much smaller initial cost than PSIR, it is clear that FrodoPIR outperforms the alternatives.

In the online phase, the client download (Figure 6 (4)) in FrodoPIR is smallest for all captured DB sizes. The server response growth rate, even for  $|\text{DB}| = 2^{20}$ , is  $< 3.6\times$ , which is significantly smaller than that of SOnionPIR ( $128\times$ ) and PSIR ( $320\times$ ). The major downside of the FrodoPIR approach is that the client query in the online phase (Figure 6 (5)) grows linearly in the size of DB, and is much larger than both SOnionPIR and PSIR — reaching 4MB for client queries when  $|\text{DB}| = 2^{20}$ . As noted in Section 5.4, this cost can be reduced using database sharding with the additional benefit of reducing server computation times, but at the cost of increasing client download sizes during the offline phase.

*Financial costs.* In the offline phase (Figure 6 (6)), PSIR provides by far the most expensive option, due to the high network bandwidth usage. The costs of SOnionPIR scale with the number of client queries. The costs of FrodoPIR include a client-independent preprocessing phase, and much lower bandwidth usage than PSIR. Therefore, for large multi-client deployments, the costs of FrodoPIR will clearly be much cheaper than both prior schemes.

The online financial costs (Figure 6 (7)) for all protocols are much smaller than in the offline phase. By far, PSIR is the most expensive protocol to run in



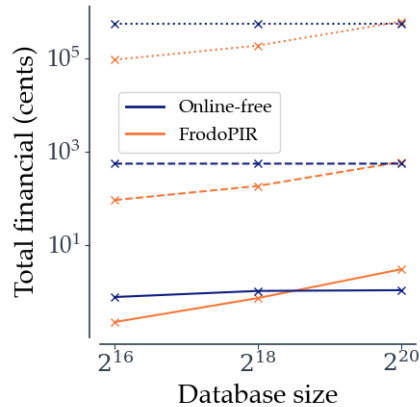
**Fig. 6.** Comparison of per-client computational, communication, and financial costs for the server when running FrodoPIR, SOnionPIR, and PSIR, assuming that each client makes  $c = 500$  queries. We include amortized costs according to various numbers of clients  $C$ , to indicate the global amortization potential of FrodoPIR. Individual charts: **(1)** Server offline computation (secs) including amortization potential over  $C$  for FrodoPIR; **(2)**: Server online computation (ms), amortized according to number of DB entries returned; **(3)**: Client offline download (KB); **(4)**: Client online download (KB); **(5)**: Client online query (KB); **(6)**: Server offline financial cost (US cents), compared for different values of  $C$ ; **(7)**: Server online financial cost (US cents).

the online phase (again, due to the high communication overhead). The costs of FrodoPIR outperform SOnionPIR, demonstrating that the trade-off between computation and communication in FrodoPIR is concretely cheaper to realize on the server-side.

*Comparison with online-free PIR.* The CHKPIR scheme [30] demonstrates entirely sublinear (amortized) running times and communication costs. However, this depends on each client launching a fairly large number of queries themselves (e.g.  $\sqrt{m}$ ). As is noted in [30], the offline phase can be implemented using the methods of PSIR or SOnionPIR, and, regardless, it is still non-amortizable across multiple clients.

To illustrate the bottleneck that the offline phase introduces from a financial perspective, we consider a PIR scheme that has *zero* online costs (which is clearly a significant underestimate for CHKPIR), and has the offline costs of SOnionPIR (sublinear in  $m$  and generally lower than PSIR). As shown in Figure 7, FrodoPIR





**Fig. 7.** Comparison of total server financial costs in FrodoPIR with an online-free PIR scheme that implements the offline phase using SOnionPIR. The costs are compared for:  $C = 1$  (solid line),  $C = 1000$  (dashed line), and  $C = 1$  million (dotted line) clients; where each client makes  $c = 500$  queries.

is cheaper to run for databases of size  $\leq 2^{18}$ , for  $c \cdot C$  queries. The costs are almost identical when  $|\text{DB}| = 2^{20}$ . We can conclude that these results, coupled with the benefits of a simple and available implementation, make FrodoPIR a very attractive option for implementing fast and scalable PIR for large multi-client systems.

### 6.3 Stateless PIR and Larger Database Elements

As mentioned in Section 2, stateless schemes tend to be less efficient than stateful schemes. However, the very recent and notably efficient Spiral PIR scheme of Menon and Wu [59] has been shown to produce low running costs across databases of various sizes and shapes. Spiral demonstrates highly promising performance for both standard PIR use-cases, and those that involve streaming large files.

To illustrate how FrodoPIR performs for different database shapes, Table 7 compares the online computational and bandwidth costs of both FrodoPIR and Spiral for each of the database types considered in [59], using the same AWS `c5n.2xlarge` EC2 instance that is used in the original work.<sup>16</sup> In the table, the *rate* of both schemes is the ratio of the response size to the retrieved database element, and the *throughput* is the ratio of the database size to the server’s computation time. In both cases, higher values are preferable.

<sup>16</sup> Note that we specifically compare with the single query variant of Spiral (with packing optimizations [59]), rather than the SpiralStream variant that is optimized for streaming use-cases.

**Table 7.** Comparison of single-threaded server-side performance between FrodoPIR and Spiral across a variety of database sizes on a `c5n.2xlarge` AWS EC2 instance. <sup>†</sup>For the database containing  $2^{18} \times 30\text{KB}$  elements, we extrapolate FrodoPIR compute performance for a database containing  $2^{16}$  elements instead, since the EC2 instance does not have sufficient memory for storing the preprocessed database of size  $2^{18}$ .

Database	Metric	Spiral	FrodoPIR
$2^{20} \times 256\text{B}$	One-time preprocessing	—	327s
	Per-client download	14MB	1.54MB
	Query size	14KB	4MB
	Resp. size	20KB	912B
	Computation	1.37s	0.16s
	Rate	0.0125	0.28
	Throughput	196MB/s	1.56GB/s
$2^{18} \times 30\text{KB}^{\dagger}$	One-time preprocessing	—	7703s
	Per-client download	18MB	166MB
	Query size	14KB	1MB
	Resp. size	86KB	96KB
	Computation	17.69s	4.27s
	Rate	0.3488	0.3125
	Throughput	434MB/s	1.76GB/s
$2^{14} \times 100\text{KB}$	One-time preprocessing	—	1605s
	Per-client download	47MB	554MB
	Query size	14KB	64KB
	Resp. size	188KB	320KB
	Computation	4.58s	0.89s
	Rate	0.5307	0.3125
	Throughput	358MB/s	1.76GB/s

Table 7 illustrates that FrodoPIR is very efficient for *narrow* databases, where each data element is relatively small, outperforming Spiral in almost all criteria. Spiral generally outperforms FrodoPIR in cases where database elements are larger: demonstrating smaller bandwidth usage and higher rate. However, FrodoPIR demonstrates faster computation across all database sizes than Spiral, and this leads to significantly higher throughput in all cases. The offline preprocessing phase of FrodoPIR can be expensive, but recall that this amortized across all clients (and queries) in the global system, and so these costs amortize away fairly quickly.

Overall, in cases where database elements are relatively small, or where fast computation is required, FrodoPIR excels compared to the recent state-of-the-art in PIR design. In cases, where database elements are large, and bandwidth and/or client storage are constrained, Spiral excels.

## 7 Discussion

*Supporting Keyword Queries.* In the interest of supporting more realistic database queries, Chor et al. constructed a PIR-by-keyword framework, where the server DB is a key-value store and client queries are keywords that recover the associated values [27]. Their framework runs multiple instances of index-based PIR as a black-box; FrodoPIR is compatible with this approach. The work of Boyle

et al. [18], based upon multi-server distributed point functions, includes direct support for keyword queries directly, but it does not appear to generalize to other PIR schemes.

As well as generic frameworks, FrodoPIR is compatible with external mechanisms for deciding keyword-to-index association. Such mechanisms include the approach detailed by Kogan and Corrigan-Gibbs [54], that furnishes the client with  $O(m)$  hash prefixes of each keyword, and associates each with a server DB index. This allows the client to learn the index that they need to query, without running multiple instances of the PIR scheme. It requires sending  $O(m)$  data to the client but which, in practice, is a very small fraction of the real database. We discuss practical costs in Appendix B.

*Optimizations for server computation.* We avoided discussing computational optimizations in this work, in favor of maintaining simplicity and configurability of FrodoPIR. However, asymptotic overheads for computing matrix multiplications have seen a variety of improvements over the last half a century [71, 29, 74, 4]. Such findings have been used in previous PIR schemes to reduce computational workloads [57, 41]. The server offline phase in FrodoPIR involves a large matrix multiplication with dimensions  $n \times m$  and  $m \times \omega$ , which would clearly benefit from sub-cubic multiplication methods. The client offline phase, involves preprocessing  $c$  queries, each involving a vector-matrix multiplication, which could be batched together into a single matrix multiplication. Furthermore, the server online phase involves a vector-matrix multiplication, for every client query. This optimization can be used by batching a number of queries together. As is observed by Lueks and Goldberg [57], this enables the server’s work to scale sublinearly in the number of client queries.

## 7.1 Applications of PIR

In Appendix B, we illustrate how efficient FrodoPIR can be when applied to the real-world of the SafeBrowsing API [44]. We list various other applications below that could also benefit. Valuable future work would identify whether FrodoPIR is a practical candidate for solving such applications.

*Certificate auditing.* Certificate Transparency (CT) is a system created to increase visibility of issued certificates. This system allows detection of misissued certificates or other forms of Certificate Authorities (CA) misbehavior, via interaction with one or more public logs. Clients should check that certificates are indeed included in these logs, but this leads to a potential privacy violation as it means that, over time, the client presents the browsing history of the user. One can apply FrodoPIR to check whether the promise of inclusion is fulfilled. Similar applications of PIR have been highlighted in concurrent work [50].

*Certificate revocation checks.* Certificate revocation checks typically use the Online Certificate Status Protocol (OCSP). This mechanism allows CAs to inform

clients if a certificate is revoked by having them query an endpoint. This mechanism, however, can violate privacy as the certificates are revealed to the CA. An alternative is to have clients download certificate revocation lists (CRLs) from endpoints maintained by CAs. This, though, comes with a huge storage overhead and the need for regular updates. FrodoPIR could be used to perform OCSP queries in a privacy-preserving manner.

*PIR for streaming.* PIR schemes such as Popcorn [47] and Spiral [59] identify PIR as a potential solution for private streaming use-cases, where clients gradually consume chunks of a large data element (such as a video). The capability of FrodoPIR for sharding the server database (Section 5.4) could make it a viable candidate in this setting.

## 7.2 Concurrent work

Concurrently to our work, the authors of [50] presented two new PIR schemes, known as SimplePIR and DoublePIR. The SimplePIR approach uses the same underlying mechanism as FrodoPIR, i.e. using LWE-based encryption reminiscent of Regev [66], and structures the database in a similar way. The main differences between the two approaches are described below.

- The database matrix is assumed to be square in [50], which means that the client upload/download are balanced.
- The parameter choices for the LWE instance are slightly different. The SimplePIR scheme uses an LWE dimension of  $n = 1024$ , and also considers moduli of size  $2^{16}$  and  $2^{32}$ . The error distribution is also chosen to be a discrete Gaussian with standard deviation parameter  $\sigma = 6.4$ . Moreover, the plaintext modulus  $\rho$  ranges from 247 to 991.

Overall, these variations lead to some differences in the perceived performance. For example, SimplePIR requires downloading a much bigger set of public parameters than in FrodoPIR: 124MB for a 1GB database, rather than  $\sim 6$ MB. The overall online communication is, however, smaller in SimplePIR, due to the client upload being smaller than in FrodoPIR (242KB, against  $\sim 4$ MB). In summary, we view these two schemes as complimentary, and further evidence to highlight the inherent configurability of using the LWE-based approach when building simple and efficient PIR schemes.

The DoublePIR approach of [50] presents an optimization that we do not consider here that refers to running SimplePIR twice: once on the original database to recover the public parameters matrix  $\mathbf{M}$ , and then again on this matrix  $\mathbf{M}$ . The observation arises from the original paper of Kushilevitz and Ostrowsky [55], that observes that only a small part of this matrix is required when computing the output of an online query (i.e. only the  $i^{\text{th}}$  row of  $\mathbf{M}$ ). This allows reducing the size of the client public parameters download to around 16MB, at the cost of increasing the amount of work that the server does during the online phase. The philosophy behind the DoublePIR optimisation is also applicable to FrodoPIR, though we do not present explicit experimental results to see if it can result in performance improvements here.

## 8 Conclusion

In this work, we built FrodoPIR. Via a simple proof-of-concept Rust implementation,<sup>17</sup> we illustrated that FrodoPIR is concretely cheaper than the previous state-of-the-art in building stateful PIR schemes, especially in large multi-client deployments. Overall, we believe that FrodoPIR is the first single-server PIR scheme that is both flexible and efficient enough to be deployed at scale, for a variety of applications.

## Acknowledgements

The authors would like to thank Benjamin Livshits, Hamed Haddadi, Joe Rowell, Muhammad Haris Mughees, Martin Albrecht, Ling Ren, Alexandra Henzinger, Matthew M. Hong, Henry Corrigan-Gibbs, and Sarah Meiklejohn for providing helpful feedback during various stages of this work. This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

## References

1. C. Aguilar Melchor, J. Barrier, L. Fousse, and M.-O. Killijian. XPIR: Private information retrieval for everyone. *PoPETs*, 2016(2):155–174, Apr. 2016.
2. M. R. Albrecht, R. Player, and S. Scott. On the concrete hardness of learning with errors. *J. Math. Cryptol.*, 9(3):169–203, 2015. (Estimates calculated Dec 2021).
3. A. Ali, T. Lepoint, S. Patel, M. Raykova, P. Schoppmann, K. Seth, and K. Yeo. Communication–Computation trade-offs in PIR. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1811–1828. USENIX Association, Aug. 2021.
4. J. Alman and V. V. Williams. A refined laser method and faster matrix multiplication. In D. Marx, editor, *32nd SODA*, pages 522–539. ACM-SIAM, Jan. 2021.
5. A. Ambainis. Upper bound on communication complexity of private information retrieval. In P. Degano, R. Gorrieri, and A. Marchetti-Spaccamela, editors, *ICALP 97*, volume 1256 of *LNCS*, pages 401–407. Springer, Heidelberg, July 1997.
6. S. Angel, H. Chen, K. Laine, and S. T. V. Setty. PIR with compressed queries and amortized query processing. In *2018 IEEE Symposium on Security and Privacy*, pages 962–979. IEEE Computer Society Press, May 2018.
7. S. Angel and S. Setty. Unobservable communication over fully untrusted infrastructure. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI’16*, page 551–569, USA, 2016. USENIX Association.
8. S. Angel and S. Setty. Unobservable communication over fully untrusted infrastructure. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 551–569, Savannah, GA, Nov. 2016. USENIX Association.
9. R. Beigel, L. Fortnow, and W. Gasarch. A tight lower bound for restricted pir protocols. *Computational Complexity*, 15:82–91, 05 2006.

<sup>17</sup> <https://github.com/brave-experiments/frodo-pir>

10. A. Beimel and Y. Ishai. Information-theoretic private information retrieval: A unified construction. In F. Orejas, P. G. Spirakis, and J. van Leeuwen, editors, *ICALP 2001*, volume 2076 of *LNCS*, pages 912–926. Springer, Heidelberg, July 2001.
11. A. Beimel, Y. Ishai, E. Kushilevitz, and I. Orlov. Share conversion and private information retrieval. In *2012 IEEE 27th Conference on Computational Complexity*, pages 258–268, 2012.
12. A. Beimel, Y. Ishai, E. Kushilevitz, and J.-F. Raymond. Breaking the  $O(n^{1/(2k-1)})$  barrier for information-theoretic private information retrieval. In *43rd FOCS*, pages 261–270. IEEE Computer Society Press, Nov. 2002.
13. A. Beimel, Y. Ishai, and T. Malkin. Reducing the servers computation in private information retrieval: PIR with preprocessing. In M. Bellare, editor, *CRYPTO 2000*, volume 1880 of *LNCS*, pages 55–73. Springer, Heidelberg, Aug. 2000.
14. S. Bell and P. Komisarczuk. An analysis of phishing blacklists: Google safe browsing, openphish, and phishtank. In *Proceedings of the Australasian Computer Science Week Multiconference, ACSW '20*, New York, NY, USA, 2020. Association for Computing Machinery.
15. D. J. Bernstein, C. Chuengsatiansup, T. Lange, and C. van Vredendaal. NTRU prime: Reducing attack surface at low cost. In C. Adams and J. Camenisch, editors, *SAC 2017*, volume 10719 of *LNCS*, pages 235–260. Springer, Heidelberg, Aug. 2017.
16. N. Borisov, G. Danezis, and I. Goldberg. DP5: A private presence service. *PoPETs*, 2015(2):4–24, Apr. 2015.
17. J. W. Bos, C. Costello, L. Ducas, I. Mironov, M. Naehrig, V. Nikolaenko, A. Raghunathan, and D. Stebila. Frodo: Take off the ring! Practical, quantum-secure key exchange from LWE. In Weippl et al. [73], pages 1006–1018.
18. E. Boyle, N. Gilboa, and Y. Ishai. Function secret sharing: Improvements and extensions. In Weippl et al. [73], pages 1292–1303.
19. E. Boyle, Y. Ishai, R. Pass, and M. Wootters. Can we access a database both locally and privately? In Kalai and Reyzin [53], pages 662–693.
20. Z. Brakerski, A. Langlois, C. Peikert, O. Regev, and D. Stehlé. Classical hardness of learning with errors. In D. Boneh, T. Roughgarden, and J. Feigenbaum, editors, *45th ACM STOC*, pages 575–584. ACM Press, June 2013.
21. C. Cachin, S. Micali, and M. Stadler. Computationally private information retrieval with polylogarithmic communication. In J. Stern, editor, *EUROCRYPT'99*, volume 1592 of *LNCS*, pages 402–414. Springer, Heidelberg, May 1999.
22. L. Caires, G. F. Italiano, L. Monteiro, C. Palamidessi, and M. Yung, editors. *ICALP 2005*, volume 3580 of *LNCS*. Springer, Heidelberg, July 2005.
23. R. Canetti, J. Holmgren, and S. Richelson. Towards doubly efficient private information retrieval. In Kalai and Reyzin [53], pages 694–726.
24. Y.-C. Chang. Single database private information retrieval with logarithmic communication. In H. Wang, J. Pieprzyk, and V. Varadharajan, editors, *ACISP 04*, volume 3108 of *LNCS*, pages 50–61. Springer, Heidelberg, July 2004.
25. H. Chen, I. Chillotti, and L. Ren. Onion ring ORAM: Efficient constant bandwidth oblivious RAM from (leveled) TFHE. In L. Cavallaro, J. Kinder, X. Wang, and J. Katz, editors, *ACM CCS 2019*, pages 345–360. ACM Press, Nov. 2019.
26. B. Chor and N. Gilboa. Computationally private information retrieval (extended abstract). In *29th ACM STOC*, pages 304–313. ACM Press, May 1997.
27. B. Chor, N. Gilboa, and M. Naor. Private information retrieval by keywords. Cryptology ePrint Archive, Report 1998/003, 1998. <https://eprint.iacr.org/1998/003>.

28. B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan. Private information retrieval. In *36th FOCS*, pages 41–50. IEEE Computer Society Press, Oct. 1995.
29. D. Coppersmith and S. Winograd. Matrix multiplication via arithmetic progressions. *J. Symb. Comput.*, 9(3):251–280, mar 1990.
30. H. Corrigan-Gibbs, A. Henzinger, and D. Kogan. Single-server private information retrieval with sublinear amortized time. In O. Dunkelman and S. Dziembowski, editors, *Advances in Cryptology – EUROCRYPT 2022*, pages 3–33, Cham, 2022. Springer International Publishing.
31. H. Corrigan-Gibbs and D. Kogan. Private information retrieval with sublinear online time. In A. Canteaut and Y. Ishai, editors, *EUROCRYPT 2020, Part I*, volume 12105 of *LNCS*, pages 44–75. Springer, Heidelberg, May 2020.
32. S. Devadas, M. van Dijk, C. W. Fletcher, L. Ren, E. Shi, and D. Wichs. Onion ORAM: A constant bandwidth blowup oblivious RAM. In E. Kushilevitz and T. Malkin, editors, *TCC 2016-A, Part II*, volume 9563 of *LNCS*, pages 145–174. Springer, Heidelberg, Jan. 2016.
33. C. Dong and L. Chen. A fast single server private information retrieval protocol with low communication cost. In M. Kutylowski and J. Vaidya, editors, *ESORICS 2014, Part I*, volume 8712 of *LNCS*, pages 380–399. Springer, Heidelberg, Sept. 2014.
34. L. Ducas, A. Durmus, T. Lepoint, and V. Lyubashevsky. Lattice signatures and bimodal Gaussians. In R. Canetti and J. A. Garay, editors, *CRYPTO 2013, Part I*, volume 8042 of *LNCS*, pages 40–56. Springer, Heidelberg, Aug. 2013.
35. Z. Dvir and S. Gopi. 2-server PIR with subpolynomial communication. *J. ACM*, 63(4):39:1–39:15, 2016.
36. K. Efremenko. 3-query locally decodable codes of subexponential length. *SIAM J. Comput.*, 41(6):1694–1703, 2012.
37. E. Fung, G. Kellaris, and D. Papadias. Combining differential privacy and PIR for efficient strong location privacy. In C. Claramunt, M. Schneider, R. C. Wong, L. Xiong, W. Loh, C. Shahabi, and K. Li, editors, *Advances in Spatial and Temporal Databases - 14th International Symposium, SSTD 2015, Hong Kong, China, August 26-28, 2015. Proceedings*, volume 9239 of *Lecture Notes in Computer Science*, pages 295–312. Springer, 2015.
38. C. Gentry and Z. Ramzan. Single-database private information retrieval with constant communication rate. In Caires et al. [22], pages 803–815.
39. T. Gerbet, A. Kumar, and C. Lauradoux. A privacy analysis of google and yandex safe browsing. In *2016 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pages 347–358, 2016.
40. N. Gilboa and Y. Ishai. Distributed point functions and their applications. In P. Q. Nguyen and E. Oswald, editors, *EUROCRYPT 2014*, volume 8441 of *LNCS*, pages 640–658. Springer, Heidelberg, May 2014.
41. I. Goldberg. Improving the robustness of private information retrieval. In *2007 IEEE Symposium on Security and Privacy*, pages 131–148. IEEE Computer Society Press, May 2007.
42. O. Goldreich. Towards a theory of software protection and simulation by oblivious RAMs. In A. Aho, editor, *19th ACM STOC*, pages 182–194. ACM Press, May 1987.
43. O. Goldreich and R. Ostrovsky. Software protection and simulation on oblivious rams. *J. ACM*, 43(3):431–473, 1996.
44. Google. SafeBrowsing API, 2008. URL: <https://safebrowsing.google.com/>. Accessed Jan. 25, 2022.



45. M. Green, W. Ladd, and I. Miers. A protocol for privately reporting ad impressions at scale. In Weippl et al. [73], pages 1591–1601.
46. T. Güneysu, V. Lyubashevsky, and T. Pöppelmann. Practical lattice-based cryptography: A signature scheme for embedded systems. In E. Prouff and P. Schumacher, editors, *CHES 2012*, volume 7428 of *LNCS*, pages 530–547. Springer, Heidelberg, Sept. 2012.
47. T. Gupta, N. Crooks, W. Mulhern, S. Setty, L. Alvisi, and M. Walfish. Scalable and private media consumption with popcorn. In *Proceedings of the 13th Usenix Conference on Networked Systems Design and Implementation*, NSDI'16, page 91–107, USA, 2016. USENIX Association.
48. A. Hamlin, R. Ostrovsky, M. Weiss, and D. Wichs. Private anonymous data access. In Y. Ishai and V. Rijmen, editors, *EUROCRYPT 2019, Part II*, volume 11477 of *LNCS*, pages 244–273. Springer, Heidelberg, May 2019.
49. R. Henry. Polynomial batch codes for efficient IT-PIR. *PoPETs*, 2016(4):202–218, Oct. 2016.
50. A. Henzinger, M. M. Hong, H. Corrigan-Gibbs, S. Meiklejohn, and V. Vaikuntanathan. One server for the price of two: Simple and fast single-server private information retrieval. Cryptology ePrint Archive, Paper 2022/949, 2022. <https://eprint.iacr.org/2022/949>.
51. A. Hülsing, J. Rijneveld, J. M. Schanck, and P. Schwabe. High-speed key encapsulation from NTRU. In W. Fischer and N. Homma, editors, *CHES 2017*, volume 10529 of *LNCS*, pages 232–252. Springer, Heidelberg, Sept. 2017.
52. Y. Ishai, E. Kushilevitz, R. Ostrovsky, and A. Sahai. Batch codes and their applications. In L. Babai, editor, *36th ACM STOC*, pages 262–271. ACM Press, June 2004.
53. Y. Kalai and L. Reyzin, editors. *TCC 2017, Part II*, volume 10678 of *LNCS*. Springer, Heidelberg, Nov. 2017.
54. D. Kogan and H. Corrigan-Gibbs. Private blocklist lookups with checklist. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 875–892. USENIX Association, Aug. 2021.
55. E. Kushilevitz and R. Ostrovsky. Replication is NOT needed: SINGLE database, computationally-private information retrieval. In *38th FOCS*, pages 364–373. IEEE Computer Society Press, Oct. 1997.
56. H. Lipmaa. An oblivious transfer protocol with log-squared communication. In J. Zhou, J. Lopez, R. H. Deng, and F. Bao, editors, *ISC 2005*, volume 3650 of *LNCS*, pages 314–328. Springer, Heidelberg, Sept. 2005.
57. W. Lueks and I. Goldberg. Sublinear scaling for multi-client private information retrieval. In R. Böhme and T. Okamoto, editors, *FC 2015*, volume 8975 of *LNCS*, pages 168–186. Springer, Heidelberg, Jan. 2015.
58. Y. Ma, K. Zhong, T. Rabin, and S. Angel. Incremental offline/online PIR (extended version). *IACR Cryptol. ePrint Arch.*, page 1438, 2021.
59. S. J. Menon and D. J. Wu. SPIRAL: fast, high-rate single-server PIR via FHE composition. In *43rd IEEE Symposium on Security and Privacy, SP 2022, San Francisco, CA, USA, May 22-26, 2022*, pages 930–947. IEEE, 2022.
60. D. Micciancio and C. Peikert. Trapdoors for lattices: Simpler, tighter, faster, smaller. In D. Pointcheval and T. Johansson, editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 700–718. Springer, Heidelberg, Apr. 2012.
61. P. Mittal, F. G. Olumofin, C. Troncoso, N. Borisov, and I. Goldberg. PIR-tor: Scalable anonymous communication using private information retrieval. In *USENIX Security 2011*. USENIX Association, Aug. 2011.



62. M. H. Mughees, H. Chen, and L. Ren. Onionpir: Response efficient single-server pir. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21*, page 2292–2306, New York, NY, USA, 2021. Association for Computing Machinery.
63. M. H. Mughees, G. Pestana, A. Davidson, and B. Livshits. Privatefetch: Scalable catalog delivery in privacy-preserving advertising. *CoRR*, abs/2109.08189, 2021.
64. J. Park and M. Tibouchi. SHECS-PIR: Somewhat homomorphic encryption-based compact and scalable private information retrieval. In L. Chen, N. Li, K. Liang, and S. A. Schneider, editors, *ESORICS 2020, Part II*, volume 12309 of *LNCS*, pages 86–106. Springer, Heidelberg, Sept. 2020.
65. S. Patel, G. Persiano, and K. Yeo. Private stateful information retrieval. In D. Lie, M. Mannan, M. Backes, and X. Wang, editors, *ACM CCS 2018*, pages 1002–1019. ACM Press, Oct. 2018.
66. O. Regev. On lattices, learning with errors, random linear codes, and cryptography. In H. N. Gabow and R. Fagin, editors, *37th ACM STOC*, pages 84–93. ACM Press, May 2005.
67. L. Ren, C. W. Fletcher, A. Kwon, E. Stefanov, E. Shi, M. van Dijk, and S. Devadas. Constants count: Practical improvements to oblivious RAM. In J. Jung and T. Holz, editors, *USENIX Security 2015*, pages 415–430. USENIX Association, Aug. 2015.
68. S. Servan-Schreiber, K. Hogan, and S. Devadas. Adveil: A private targeted-advertising ecosystem. Cryptology ePrint Archive, Report 2021/1032, 2021. <https://ia.cr/2021/1032>.
69. R. Sion and B. Carbunar. On the practicality of private information retrieval. In *NDSS 2007*. The Internet Society, Feb. / Mar. 2007.
70. E. Stefanov, M. van Dijk, E. Shi, T. H. Chan, C. W. Fletcher, L. Ren, X. Yu, and S. Devadas. Path ORAM: an extremely simple oblivious RAM protocol. *J. ACM*, 65(4):18:1–18:26, 2018.
71. V. Strassen. Gaussian elimination is not optimal. *Numer. Math.*, 13(4):354–356, aug 1969.
72. S. Wehner and R. Wolf. Improved lower bounds for locally decodable codes and private information retrieval. In Caires et al. [22], pages 1424–1436.
73. E. R. Weippl, S. Katzenbeisser, C. Kruegel, A. C. Myers, and S. Halevi, editors. *ACM CCS 2016*. ACM Press, Oct. 2016.
74. V. V. Williams. Multiplying matrices faster than coppersmith-winograd. In H. J. Karloff and T. Pitassi, editors, *44th ACM STOC*, pages 887–898. ACM Press, May 2012.
75. S. Yekhanin. Towards 3-query locally decodable codes of subexponential length. *J. ACM*, 55(1):1:1–1:16, 2008.

## A Proofs from Section 4

### A.1 Proof of Theorem 2

Let  $\tilde{\mathbf{b}} \leftarrow \mathbf{b} + \mathbf{f}_i$ , where  $i$  is the requested index of DB by the client. As laid out in Section 4.4, we must show that Equation (6) holds, with all but negligible probability. Firstly, note that since  $\mathbf{e} \leftarrow (\chi)^m$ , then by Corollary 2, we have that  $\|\mathbf{e} \cdot \mathbf{D}\|_\infty \leq 4\rho\sqrt{m}$  with high probability. This follows because the number

of samples  $m$  is very large and by assuming that each entry in  $\mathbf{D}$  is equal to  $\rho = \|\mathbf{D}\|_\infty$ . Consequently:

$$\begin{aligned} \lfloor (\mathbf{e} + \mathbf{f}_i)^T \cdot \mathbf{D} \rfloor_\rho &= \lfloor \rho/q \cdot (\mathbf{e}^T \cdot \mathbf{D} + \mathbf{f}_i^T \cdot \mathbf{D}) \rfloor \\ &= \lfloor \rho/q \cdot (\mathbf{e}^T \cdot \mathbf{D}) + \mathbf{D}[i] \rfloor \\ &= \lfloor \mathbf{y} + \mathbf{D}[i] \rfloor, \end{aligned} \tag{9}$$

where  $\mathbf{y} = \rho/q \cdot (\mathbf{e}^T \cdot \mathbf{D})$  and  $\mathbf{D}[i] \in \mathbb{Z}_q^\omega$  is the  $i^{\text{th}}$  row of  $\mathbf{D}$  (interpreted as a column vector). Therefore,  $\|\mathbf{y}\|_\infty < 4\rho^2\sqrt{m}/q = 1/2$  by the statement of the theorem and, as a consequence:

$$\lfloor (\mathbf{e} + \mathbf{f}_i)^T \cdot \mathbf{D} \rfloor_\rho = \mathbf{D}[i], \tag{10}$$

which is the correct output of the protocol.  $\square$

## A.2 Proof of Theorem 3

Let  $\text{ip} \leftarrow \text{FPIR.ssetup}(1^\lambda)$ ,  $\text{pp} \leftarrow \text{FPIR.spreproc}(\text{DB})$ ,  $\text{st} \leftarrow \text{FPIR.cpreproc}(\text{pp})$ , let  $i_0, i_1 \leftarrow \mathcal{A}(\text{ip}, \text{pp})$ , let  $b \leftarrow_{\mathcal{S}} \{0, 1\}$ , and let  $\tilde{\mathbf{b}}_b \leftarrow \text{FPIR.query}(\text{st}, i_b)$ . In particular, we have that  $\tilde{\mathbf{b}}_b = \mathbf{s}^T \cdot \mathbf{A} + \mathbf{e}^T + \mathbf{f}_{i_b}^T \in \mathbb{Z}_q^m$ , for  $\mathbf{A} \in \text{st}$ ,  $\mathbf{s} \leftarrow (\chi)^n$ ,  $\mathbf{e} \leftarrow (\chi)^m$ ,  $\mathbf{A} \leftarrow \text{PRG}(\mu, n, m, q)$ , and  $\mathbf{f}_{i_b}$  the  $m$ -dimensional vector of all zeroes except where  $\mathbf{f}_{i_b}[i_b] = q/\rho$ . Clearly, we can show that FPIR is secure if we can show that the output of  $\text{FPIR.query}$  is distributed uniformly.

Firstly, note that  $\mathbf{A}$  is sampled as the output of a pseudorandom generator, therefore, it is indistinguishable from  $\mathbf{A} \leftarrow_{\mathcal{S}} \mathbb{Z}_q^{n \times m}$ . Let  $\mathcal{A}$  be an adversary in the  $\text{LWE}_{q,n,m,\chi}$  decisional security game (Definition 1), receiving  $(\mathbf{A}, \mathbf{u})$  as the challenge, and let  $\mathcal{S}$  be an adversary in the PIR 1-query indistinguishability game (Definition 4). When  $\mathcal{A}$  receives the sample in Equation (1),  $\mathbf{b}$  and  $\tilde{\mathbf{b}}$  are distributed identically, and when it receives the sample in Equation (2), then  $\mathbf{b} \leftarrow_{\mathcal{S}} \mathbb{Z}_q^m$ . Therefore, the adversary  $\mathcal{A}$  can simulate the client query to  $\mathcal{S}$  by simply sending  $\tilde{\mathbf{b}} = \mathbf{u} + \mathbf{f}_{i_b}$  for  $b \leftarrow_{\mathcal{S}} \{0, 1\}$ . When  $\mathcal{S}$  returns their guess  $b' \in \{0, 1\}$  to  $\mathcal{A}$ ,  $\mathcal{A}$  checks if  $b' \stackrel{?}{=} b$ .

Clearly, whatever advantage  $\epsilon$  that  $\mathcal{S}$  has in guessing the correct value of  $b$ , immediately translates to  $\mathcal{A}$  having advantage  $\epsilon$  in the decisional  $\text{LWE}_{q,n,m,\chi}$  security game. Since we assume that  $\text{LWE}_{q,n,m,\chi}$  is difficult to solve, we therefore conclude that  $\epsilon \leq \text{negl}(\lambda)$ .

To conclude, in the case that  $\mathbf{b}$  is sampled uniformly, then the adversary has no advantage in distinguishing since  $\tilde{\mathbf{b}}$  is distributed uniformly.  $\square$

### A.3 Proof of Corollary 3

We can construct a matrix  $\tilde{\mathbf{B}}$  from each query  $\tilde{\mathbf{b}}_j$  ( $j \in [\ell]$ ) that  $\mathcal{S}$  observes with the following structure:

$$\begin{aligned}
 \tilde{\mathbf{B}} &= \left[ \tilde{\mathbf{b}}_1 \parallel \cdots \parallel \tilde{\mathbf{b}}_\ell \right] \\
 &= \left[ (\mathbf{s}_1^T \cdot \mathbf{A} + \mathbf{e}_1^T)^T + \mathbf{f}_{i_1} \parallel \cdots \parallel (\mathbf{s}_\ell^T \cdot \mathbf{A} + \mathbf{e}_\ell^T)^T + \mathbf{f}_{i_\ell} \right] \\
 &= (\mathbf{s}_1 \parallel \cdots \parallel \mathbf{s}_\ell)^T \cdot \mathbf{A} + [\mathbf{e}_1 \parallel \cdots \parallel \mathbf{e}_\ell]^T + [\mathbf{f}_{i_1} \parallel \cdots \parallel \mathbf{f}_{i_\ell}] \\
 &= \mathbf{S} \cdot \mathbf{A} + \mathbf{E} + \mathbf{F} \in \mathbb{Z}_q^{\ell \times m}.
 \end{aligned} \tag{11}$$

Using the same proof strategy as in Theorem 3, we can use  $\mathcal{A}$  as an adversary attempting to decide in the decisional  $\text{MatLWE}_{q,n,m,\chi,\ell}$  security game (Definition 2). This illustrates that  $\mathcal{A}$  has advantage equal to that which  $\mathcal{S}$  has in deciding the uniformity of  $\tilde{\mathbf{B}}$ . Furthermore, by Corollary 1, we know that  $\epsilon = \ell \cdot \nu$ , where  $\nu$  is the max advantage of winning in  $\text{LWE}_{q,n,m,\chi}$ . Since  $\ell = \text{poly}(\lambda)$ , then  $\epsilon = \text{poly}(\lambda) \cdot \text{negl}(\lambda) = \text{negl}(\lambda)$ .  $\square$

## B SafeBrowsing Example

Major browsers such as Google Chrome, Firefox, and Brave integrate a security service run by Google and known as the SafeBrowsing API [44]. SafeBrowsing allows browsers to verify if online resources and webpages that the user requests are “safe”. If a resource has been flagged as “unsafe”, the user is warned by the browser and asked to explicitly consent visiting the website that contains the unsafe resource. The SafeBrowsing service relies on a list of blocked resources maintained by Google, and it exposes an API that informs the browser if a resource is part of the blocked list. The downside of serving queries to the SafeBrowsing API remotely is that clients would effectively reveal their browsing patterns to Google. It is clear that it will be important to build mechanisms that preserve client privacy from third parties (like Google, in this case), while still being able to inform users if they are about to load malicious content.

### B.1 Current SafeBrowsing Implementation

*Local storage.* In order to avoid calling the remote API for *every* resource, the entire SafeBrowsing blocklist could be shipped with each browser, but storing the full blocklist ( $> 90\text{MB}$ ) may be beyond some clients. Consequently, every browser instead stores a compressed and probabilistic data structure that contains an approximate view of the SafeBrowsing blocklist. This local data structure allows performing probabilistic checks of inclusion, with non-negligible chances of *false-positives* occurring and no chance of *false-negatives*. Due to the rate of potential *false positives*, if an inclusion check returns positive (i.e. an unsafe resource), the browser remains uncertain. To remove the uncertainty, the browser confirms

if the resource is unsafe by calling the remote SafeBrowsing API. Thus, the browser only relies on the remote API call to SafeBrowsing services when the set inclusion against the local data structure returns a potential false positive. This mechanism reduces considerably the number of remote API calls at the expense of storing a compressed, space optimized data structure locally in the browser.

The local blocklist is a set of 32-bit hashes of the resource URI, and the full SafeBrowsing blocklist consists of a key-value database mapping a 32-bit hash to a SHA256 hash of a blocked resource URI. The local blocklist results in storage and bandwidth that is  $8\times$  smaller than the full SafeBrowsing blocklist. We summarize the two distinct phases of SafeBrowsing checks in the following.

1. **(Phase 1: Local check)** First, the browser computes the 32-bit hash of the resource URI that has been requested, and checks if the 32-bit hash is part of the local storage. If the set inclusion operation returns ‘false’ (i.e. the hash of the resource does not exist in the local data structure), then the browser considers the resource safe and proceeds. If the set inclusion operation returns ‘true’ (i.e. the 32-bit resource hash is part of the local block list), the client proceeds to the next phase.
2. **(Phase 2: Remote check)** When Phase 1 identifies a *possibly* unsafe resource, the browser needs to confirm whether the resource is a false positive or not. To do so, it requests the full SHA256 hash of the resource’s URI by querying the remote SafeBrowsing API for the 32-bit hash of the resource URI computed in Phase 1. If the full SHA256 hash returned by the remote SafeBrowsing API matches with the SHA256 of the resource URI, then the resource is part of the SafeBrowsing blocklist, and the browser considers the resource unsafe.

*Privacy considerations.* The remote SafeBrowsing resource check (Phase 2) requires the browser to explicitly include the 32-bit hash identifying the resource that is being checked for inclusion on the SafeBrowsing blocklist. As noted by [14, 39], this request leaks information about the browsing history of the user, as the SafeBrowsing API service is able to learn which content a particular user is interested in. Over time, this information can be used by the SafeBrowsing service provider to construct a behavior profiling of web users without their consent.

## B.2 SafeBrowsing via FrodoPIR

FrodoPIR can be used to implement the remote SafeBrowsing API service, such that no leakage occurs during the remote SafeBrowsing API check. The intuition is that, once the index that must be queried is known to the client, the remote check can be performed via a PIR query to a remote FrodoPIR database, that stores all the SHA256 hashes of the unsafe URIs. Given the privacy guarantees of FrodoPIR, the client does not leak which resource ID is being queried.

*Requirements.* Based on the estimates provided by [54], the current SafeBrowsing blocklist contains about 3 million entries. The blocklist grows at a rate of 30,000 new entries per week. Each of the values in the database consists of a SHA256 hash of the content URI.

*Mapping URL hashes to query indices.* As in [54], we will assume that local blocklist is augmented to include the index that must be queried in the online database. That is, when the client finds a match in their local blocklist, they use the corresponding index  $i$  that is included to make a query for element  $i$  in the remote server database.

*Database configuration.* As shown in Section 5, FrodoPIR provides a high degree of flexibility, allowing developers to choose which trade-offs to make when deploying an instance of the PIR database. We now suggest the following database configuration to implement FrodoPIR for SafeBrowsing:

- We choose  $q = 2^{32}$  and  $n = 1774$ , which should be satisfactory for even the large number of clients using major Internet browsers that integrate with the SafeBrowsing API. According to [2], this provides 128-bit security for  $2^{52}$  client queries. In other words, this allows 4 billion clients to each make 1 million queries, which should be more than enough.
- We require  $w = 256$  bits for storing each URI hash in the server database.
- Let  $\tilde{m}$  be the total number of elements in the SafeBrowsing database. We require that  $\tilde{m} \geq 2^{21}$  to accommodate all the 3 million entries and subsequent updates [54]. However, we leverage sharding to break down the databases into smaller sub-databases, as explained in Section 5.4. Assuming that FrodoPIR is running on a machine with 16 cores, we can split the blocklist into  $s = 16$  sub-databases, resulting in setting  $m = 2^{18}$  per shard. This provides a database with total size  $2^{22}$ , which is enough to store the entire blocklist.
- Given  $w$ ,  $m$  and  $q$ , we set  $\rho = 2^{10}$  so that the correctness guarantee from Theorem 2 holds true.
- We calculate  $\omega = m / \log(\rho) = 26$  as described in Section 5.
- The local blocklist that each client must download contains  $32 \cdot (m + 1)$  bits to include each 32-bit hash prefix plus the corresponding 32-bit index.

We leverage sharding in two different ways. On one hand, to decrease the size of the database by splitting it into sub-databases, allowing us to reduce the size  $m$  of each sub-database, and to optimize both user and server performance and bandwidth. In addition, sharding is used to implement a low-cost database update mechanism. Updates to the blocklist happen by adding elements to one sub-database only, in turn requiring clients to derive new parameters only for a single shard at every update, as explained in Section 5.4. This is possible in SafeBrowsing because DB updates are typically only additions, and thus deletion of old content in previous shards is rarely required [54].

**Table 8.** Performance analysis of the FrodoPIR scheme when communicating with a single database shard, using the parameters defined in Section B.2.

Offline	Client download (KB)	180
	Database preprocessing (s)	28.555
	Client derive params (s)	2.2281
	Client query preprocessing (s)	0.573
Online	Client query (KB)	1024
	Server response (KB)	0.1
	Client query (ms)	0.097
	Server response (ms)	5.223
	Client output (ms)	0.012

### B.3 Implementation and Raw Costs

We set up the experimental environment, and report results in Table 8, corresponding to the raw costs of using the FrodoPIR scheme on the aforementioned parameters. We run all experiments as single-threaded processes on the same Amazon `t2.2xlarge` EC2 instance, with 8 CPU cores and 32GB of RAM, as was used in Section 6.

### B.4 Performance Analysis

From Table 8, we estimate the performance of instantiating the SafeBrowsing API for a single database shard using FrodoPIR, using the parameter set defined in Section B.2. Our extrapolations are based on the following set of usage model assumptions that are taken from the previous work of Kogan and Corrigan-Gibbs [54] on exploring usage of PIR for satisfying the demands of SafeBrowsing.

- On average, clients launch a query every 44 minutes. Assuming 12 hours of daily usage, this leads to approximately 16 queries per day.
- On average, the server database is updated every 94 minutes. This leads to around 16 DB updates per day, with a weekly addition of around 30,000 records.
- The server is a collection of  $Z$  replicas that are distributed globally, that each independently possess and process queries on the same database. Any client query can be fulfilled by a single server.
- Client storage must be, at least, a constant factor smaller than the entire SafeBrowsing database size.

*Database initialization and updates.* The main server initializes the sub-database, public parameters, and local blacklist for each individual shard. Each of these remain static for a monthly period and are downloaded by each server replica. When the main server initializes, or rotates the matrix  $\mathbf{A}$ , it posts the public parameters  $\mathbf{pp} = (\mu, M = \{\mathbf{M}_i = \mathbf{A} \cdot \mathbf{D}_i\}_{i \in [16]})$  and local blacklists to a public

location that clients can access and download from. Note that  $\mathbf{M}_i \in \mathbb{Z}_q^{m \times \omega}$  corresponds to the public parameters made available for each sub-database.

Based on our usage model, we will assume that there are 16 database updates made by the server, each containing 268 records. We assume that clients each download and process 8 updates per day. Each database update touches a single shard  $\text{DB}_i$ , and results in uploading a new value of  $\mathbf{M}_i$ .

*Client processing.* Client preprocessing amounts to preprocessing 16 queries per day, using the server provided parameters  $\mathbf{pp}$ . After every update, the client needs to regenerate the remaining preprocessed state that is associated with the sub-database that was updated. Recall that the client stores:

$$X = (\mathbf{b}_j = \mathbf{s}^T \cdot \mathbf{A} + \mathbf{e}^T, C_j = \{\mathbf{c}_i = \mathbf{s}_j^T \cdot \mathbf{M}_i\}_{i \in [16]})_{j \in [16]}$$

for each of the  $j \in [16]$  queries that the client will launch, and for each of the  $i \in [16]$  database shards. The client must also store each  $\mathbf{s}_j$  that it samples, for responding to server updates as well as the local blocklist.

Overall, at the start of each day, the client rederives  $\mathbf{A} \leftarrow \text{PRG}(\mu, n, m, q)$ , and computes the set  $X$ . Every time that the client makes a remote query it removes a pair  $(\mathbf{b}_j, C_j)$  from storage, and sends  $\tilde{\mathbf{b}}_j = \mathbf{b} + \mathbf{f}_i$  to the server, for query index  $i$  computed during the local blocklist check. Whenever the server issues a database update for shard  $i$ , the client redownloads  $\mathbf{M}_i$  and the local blocklist, and uses  $\mathbf{s}_j$  to update  $\mathbf{c}_i = \mathbf{s}_j^T \cdot \mathbf{M}_i \in C_j$ , for each remaining  $j$  (i.e. unused preprocessed query data). According to Table 8, we have the following (per-day) client computational costs.

- A single derivation of  $\mathbf{A}$ .
- preprocessing of 16 queries for each of the 16 shards.
- Updating of  $2 \sum_{a=1}^7 a = 56$  queries per day.
- 16 individual online queries.

We ignore the cost of running queries on the 32-bit hashes in the local blocklist, since these are negligible by comparison. Furthermore, the per shard cost of updating preprocessed query data is almost zero. Therefore, we calculate the total CPU costs of each client to amount to  $32.96 + 16 \cdot 0.47 + 16 * 0.00025 = 40.48$  seconds per day.

*Client download.* The initial client download of public parameters is equal to  $128 + 16 \cdot (n\omega \log(q)) = 23,615,616$  bits, which corresponds to around 2.82MB. The total size of the local blocklist is approximately  $32 \cdot 3\text{million}$  bits, which is equal to 11.44MB. The running download cost per-day is calculated as  $16\omega \log(q) + 8n\omega \log(q) + 32 * 268 = 11,829,632$  bits, which is roughly 1.41MB.

*Client query.* The client query is linear in the size of a single shard, which has a maximum of  $2^{18}$  elements. Therefore, each query is around 1MB in size, based on the costs from Table 6. As a consequence, this results in roughly 16MB of additional communication per-day.

**Table 9.** Comparison of instantiating the SafeBrowsing API using either FrodoPIR, or via the two-server PIR schemes of [54]. Estimated costs are marked with asterisks.

Performance indicators	Non-private	dpfPIR	ooPIR	FrodoPIR
Servers per 1B users	143	9047	1348	9778*
Latency (ms)	90	122	91	90*
Client init (sec)	3.1	2.6	13.3	32.96*
Client running (sec/month)	0.5	0.8	8.0	1272.0*
Initial communication (MB)	5.0	5.0	10.3	2.82
Online communication (MB/month)	3.0	3.6	9.0	539.7
Max storage (MB)	4.5	4.5	26.1	30.69*

*Client storage.* The client needs  $\sim 1\text{MB}$  to store each preprocessed query, and each secret vector  $s_j$ , for  $j \in [16]$ . In total, this represents about 16MB of required storage. Secondly, the client must store the local prefix table for the SafeBrowsing API which amounts to storing a further 11.44MB of data. Thirdly, the client stores the public parameters made available by the server, which totals 2.82MB. Overall, the maximum client storage overhead is  $\sim 30.69\text{MB}$ , which is a  $91.55/30.26 = 3.0\times$  saving compared with storing the original database. As the client makes queries, it deletes used preprocessed data, and so this storage overhead will decrease as the day progresses.

*Server processing.* The non-private SafeBrowsing API has an average latency of around 90ms per client query [54]. This is achieved using  $Z = 143$  servers answering client queries. Note that a single FrodoPIR server can answer a single client query in  $\sim 5\text{ms}$  (Table 8). We assume that 1 billion queries are received uniformly in 90ms windows over a 44 minute period.<sup>18</sup> Therefore, in each 90ms window around 29334 client queries are received. Further, we assume that each server can answer 3 client queries in 90ms (including time taken to receive and respond to the client HTTP request). To achieve this, we would need at least 9778 individual servers each answering queries on the same FrodoPIR database for servicing 1 billion clients. Clearly, this is much more expensive than running the non-private version of SafeBrowsing, but such a number of servers is still within the realms of practicality, whilst preserving client privacy.

*Comparison with [54].* The work of Kogan and Corrigan-Gibbs presents two PIR-based constructions for running the SafeBrowsing API, one based on PIR from distributed point functions (dpfPIR), and the other based on offline-online PIR (ooPIR). Both schemes require two non-colluding servers. We compare the performance of running the SafeBrowsing API using FrodoPIR against both dpfPIR and ooPIR in Table 9.

Clearly, FrodoPIR involves heavier usage costs compared to all known solutions, either non-private or using multi-server PIR. As previously highlighted, a limitation of the FrodoPIR scheme is the client request size, which makes up

<sup>18</sup> In other words, simulating 1 query from every client every 44 minutes.



**Table 10.** Comparison of FrodoPIR overheads when choosing less conservative security parameterizations.

Number of DB items ( $\log(m)$ )	16			20		
Lattice dimension ( $n$ )	1288	1572	1774	1288	1572	1774
Client download (KB)	4125.6	5035.3	5682.5	4583.5	5594.1	6313.1
Database preprocessing (s)	82.705	96.74	104.57	1439.2	1697.5	1895.2
Client derive params (s)	0.47	0.5506	0.5826	7.22	8.79	9.25
Client query preprocessing (s)	0.135	0.147	0.147	1.933	2.149	2.343

a large proportion of the total communication (496MB per month, as opposed to 43.7MB of download). The client computation is also much heavier than in multi-server PIR, due to the requirement for computing high-dimensional cryptographic operations when preprocessing queries.

Otherwise, our estimates suggest that FrodoPIR can provide adequate performance for operators where non-colluding PIR servers are impossible to set up. However, it is worth noting that the experimental analysis of [54] provides significantly more detail than we do here. Our goal is to give a broad understanding of the increased overheads of using FrodoPIR.

## C Less Conservative LWE Parameters

Our security analysis assumes that the  $\lambda$ -bit security of the Matrix LWE problem ( $\text{MatLWE}_{q,n,m,\chi,\ell}$ ) is calculated as the  $\nu$ -bit security of the underlying LWE problem, minus the logarithm of the number of queries that are launched ( $\ell$ ). This analysis is conservative for two reasons: (1) the number of queries that we protect against with our parameter choice ( $2^{52}$ ) is very large; and (2) it’s not clear that lattice cryptanalysis can exploit the Matrix LWE problem any easier than LWE. Therefore, choosing lattice parameters that are smaller may preserve adequate security, while improving efficiency.

Here, we discuss the impact on efficiency of reducing the LWE dimension,  $n$ . In particular, Section 5, we chose  $n = 1774$ . If we choose  $n = 1572$ , then this provides 128-bit security against 1 billion queries, using the same cost model for estimating the hardness of Matrix LWE. However, if we simply focus on achieving 128-bit security for the underlying Ternary LWE instance, then we can choose  $n = 1228$  instead.<sup>19</sup>

In Table 10, we highlight how costs change as the lattice dimension reduces (for databases of size  $2^{16}$  and  $2^{20}$ ), when compared with the original performance for  $n = 1774$  from Section 6. The online phase is largely unaffected by the lattice dimension, so we omit measurements for such functions.

<sup>19</sup> See Appendix D for the lattice estimation outputs that we use for calculating these dimensions.

Overall, we see that bandwidth reduces significantly — by over 11% for  $n = 1572$ , and 27% for  $n = 1288$  — which will translate into notable financial savings when servicing queries from large numbers of clients. Furthermore, client storage requirements can be reduced by the same amount. Computational workloads are also reduced but less significantly, particularly because the server preprocessing is amortized across all clients anyway. However, client derivation of  $\mathbf{A}$  sees an approximate 20% reduction, which may be notable when considering low-powered clients.

## D Lattice Estimation

For calculating the security estimates of LWE parameters, we used the lattice estimator of [2]. Specifically, we used the code available at <https://github.com/malb/lattice-estimator>, from commit: f9dc7c625d93b9c645c56bf9dfd3d4ec202f17d1. The security estimations were obtained using the following code and corresponding output.<sup>20</sup>

```

1  from estimator import *;
2
3  # n is the lattice dimension that is used
4  for n in [1288, 1572, 1774]:
5      LWE.primal_usvp(
6          LWE.Parameters(
7              n=n,
8              q=2**32,
9              Xs=ND.Uniform(-1,1),
10             Xe=ND.Uniform(-1,1),
11             m=infinity
12         )
13     )

```

```

1  rop: 2^128.2, red: 2^128.2, delta: 1.004425, beta: 343, d:
   2447, tag: usvp
2  rop: 2^158.1, red: 2^158.1, delta: 1.003668, beta: 450, d:
   2972, tag: usvp
3  rop: 2^180.0, red: 2^180.0, delta: 1.003274, beta: 528, d:
   3376, tag: usvp

```

<sup>20</sup> The output is modified slightly to omit non-utf8 characters.