

# SophOMR: Improved Oblivious Message Retrieval from SIMD-Aware Homomorphic Compression

Keewoo Lee  
UC Berkeley

Yongdong Yeo  
Seoul National University

## Abstract

Privacy-preserving blockchains and private messaging services that ensure receiver-privacy face a significant UX challenge: each client must scan every payload posted on the public bulletin board individually to avoid missing messages intended for them. Oblivious Message Retrieval (OMR) addresses this issue by securely outsourcing this expensive scanning process to a service provider using Homomorphic Encryption (HE).

In this work, we propose a new OMR scheme that substantially improves upon the previous state-of-the-art, PerfOMR (USENIX Security’24). Our implementation demonstrates reductions of 3.3x in runtime, 2.2x in digest size, and 1.3x in key size, in a scenario with 65536 payloads (each 612 bytes), of which up to 50 are pertinent.

At the core of these improvements is a new homomorphic compression mechanism, where ciphertexts of length proportional to the number of total payloads are compressed into a digest whose length is proportional to the upper bound on the number of pertinent payloads. Unlike previous approaches, our scheme fully exploits the native homomorphic SIMD structure of the underlying HE scheme, significantly enhancing efficiency. In the setting described above, our compression scheme achieves 7.4x speedup compared to PerfOMR.

## 1 Introduction

Consider a classic scenario where a sender wants to deliver a message to a receiver without revealing any information to eavesdroppers. While end-to-end encryption ensures *message-privacy*, it alone does not guarantee *metadata-privacy*. In particular, eavesdroppers may still gather details such as who is communicating, when, and how often, which can enable traffic analysis [BMS01, MD04] and compromise privacy.<sup>12</sup> As a

<sup>1</sup><https://www.propublica.org/article/how-facebook-undermines-privacy-protections-for-its-2-billion-whatsapp-users>

<sup>2</sup><https://www.forbes.com/sites/thomasbrewster/2022/02/23/meet-the-secretive-surveillance-wizards-helping-the-fbi-and-ice-wiretap-facebook-and-google-users/>

former NSA general counsel said:<sup>3</sup> “Metadata absolutely tells you everything about somebody’s life. If you have enough metadata you don’t really need content.” In this regard, mechanisms for metadata-privacy are critical in private messaging systems [WCFJ12, CBM15, vdHLZZ15, AS16] and privacy-preserving cryptocurrencies [BCG<sup>+</sup>14, Azt, Noe15].

In this work, we specifically focus on *receiver-privacy*, which aims to hide the receiver’s identity. A plausible solution to receiver-privacy is the *public bulletin board* approach. In this framework, the sender encrypts the message using the receiver’s public key and posts it on a public bulletin board (e.g., blockchain). Over time, the bulletin board accumulates messages from various senders to various receivers. To retrieve their messages, a receiver scans the entire bulletin board and attempts to decrypt each message one by one. Messages that decrypt correctly are those intended for the receiver, while unsuccessful decryptions are ignored. This framework ensures that eavesdroppers without access to the decryption key cannot link messages to their receivers.

In fact, many privacy-preserving blockchain projects that aim to provide receiver-privacy are essentially following this approach (e.g., Zcash [BCG<sup>+</sup>14, HBHW], Aztec [Azt], Monero [Noe15], ERC-5564 [WSDB22] for Stealth Address). However, this approach poses a significant UX challenge: clients must individually scan every payload posted on a public bulletin board to avoid missing messages intended for them. In the context of privacy-preserving blockchains, this means that even light clients are required to scan all transactions.

One approach to addressing this issue is to *securely outsource* the expensive scan to a powerful server. *Homomorphic Encryption* (HE) [RAD78, Gen09], a cryptosystem that allows computation on encrypted data, is a natural cryptographic primitive for such a secure outsourcing scenario [But20b, But20a]. With HE, we can delegate the scan to an untrusted server without compromising the privacy of receivers.

Liu and Tromer [LT22] formalized this HE-based secure

<sup>3</sup><https://www.wired.com/2015/03/data-and-goliath-nsa-metadata-spying-your-secrets/>

scan-outsourcing solution as a primitive called *Oblivious Message Retrieval (OMR)*. They also explored various design choices and optimization techniques and provided a proof-of-concept implementation of OMR. However, this initial scheme required over 20 hours to process  $2^{19}$  payloads (155 ms/payload), making it impractical for real-world use.

A follow-up work, PerfOMR [LTW24b] by Liu-Tromer-Wang, advanced the practicality of OMR through further optimizations, achieving significantly improved efficiency. Specifically, PerfOMR reported  $15\times$  faster runtime compared to the original OMR construction [LT22].

## 1.1 Our Contribution

In this work, we propose a new OMR scheme, SophOMR, that substantially improves upon PerfOMR [LTW24b]. Specifically, it takes less than 3 minutes to process  $2^{16}$  payloads in a single thread (2.5 ms/payload).

- At the core of the improvement is our new *homomorphic compression* mechanism based on [CLPY24]. The main advantage of our approach lies in its compatibility with the homomorphic SIMD<sup>4</sup> structure of the underlying HE scheme. This SIMD-awareness avoids complexities found in prior works [LT22, LTW24b], resulting in a more efficient and arguably simpler scheme. Notably, SophOMR requires only  $O(\sqrt{N})$  homomorphic rotations, in contrast to PerfOMR’s  $O(N)$  rotations. In our primary setting, this yields a  $7.4\times$  speedup in homomorphic compression. See Sec. 1.2 for further details.
- We also suggest several other optimization techniques to further improve the performance of the scheme: ring-switching (Sec. 4.2), hybrid key-switching (Sec. 4.3), and BSGS-style homomorphic matrix multiplication (Sec. 4.4).
- We implement SophOMR in a C++ library<sup>5</sup> and compare its concrete performance to that of PerfOMR. Our implementation demonstrates reductions of  $3.3\times$  in runtime,  $2.2\times$  in digest size, and  $1.3\times$  in key size, in a scenario with 65536 payloads (each 612 bytes), of which up to 50 are pertinent.
- We also consider OMD (Oblivious Message Detection) variants, where clients retrieve only pertinent indices rather than the actual payloads (Sec. 4.5). OMD can be combined with PIR (Private Information Retrieval) schemes [CGKS95, KO97] to achieve OMR-like functionality. While this approach introduces additional communication rounds, it can inherit beneficial features from the underlying PIR scheme. In the same setting, our OMD variant achieves an overall  $4.4\times$  speedup and a specific  $335\times$  speedup in homomorphic compression compared to PerfOMR (Sec. 5).

<sup>4</sup>Single Instruction, Multiple Data

<sup>5</sup><https://github.com/keewoolee/SophOMR>

## 1.2 Technical Overview

### 1.2.1 System Model

More precisely, OMR operates as follows (Fig. 1). For more details, refer to Sec. 2.5 and Sec. 3.

- ① Each user runs KeyGen, publishes their signal key  $pk_{sig}$ , and send the detection key  $pk_{det}$  to Detector.
- ② Sender, who wants to deliver a payload<sup>6</sup>  $x$  to Receiver, posts the pair  $(sig, x)$  on the bulletin board BB, tagging the payload with a signal  $sig$  generated using the Receiver’s signal key  $pk_{sig}$ . The signal enables only the intended Receiver to detect the payload, while others cannot.
- ③ Upon Receiver’s requests, Detector uses the Receiver’s detection key  $pk_{det}$  to summarize the bulletin board BB into a digest D and sends it to Receiver.
- ④ Finally, Receiver decodes the digest D to recover its pertinent payloads.

At a high level, Detector’s task (③) is *homomorphic detection* of the signals, given homomorphically encrypted Receiver’s secret  $sk$  as a part of the detection key  $pk_{det}$ . Through this process, Detector obtains an encrypted sparse binary vector, where the ones indicate the indices pertinent to Receiver. This encrypted binary vector is then used to mask the payloads on Bulletin Board, preserving only the pertinent payloads while zeroizing the rest. However, notice that the result is as large as the entire Bulletin Board; the Detector cannot simply discard irrelevant payloads since the detection result is encrypted. Thus, homomorphic detection alone is insufficient to achieve a *compact* digest.

This is where *homomorphic compression* comes into play. It is a technique that obviously compresses a lengthy, homomorphically encrypted sparse vector into a digest whose size depends only on an upper bound for the number of non-zero components in the input ciphertexts. By applying homomorphic compression to the output of homomorphic detection, we can condense the sparse vector of pertinent payloads into a compact digest.

**Threat Model.** The security of OMR (Sec. 2.5) considers a passive adversary attempting to link specific messages to their receivers. The adversary has access to the public bulletin board, public keys, and any communication among the parties in the system. It may also collude with any party except the sender and receiver of the target messages.

<sup>6</sup>Note that we consider payload  $x$  as just a bitstring. In applications, these payloads are typically end-to-end encrypted.

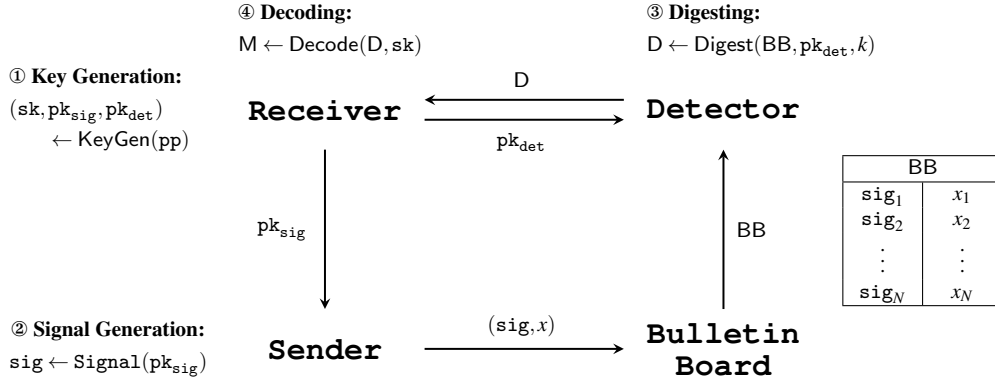


Figure 1: Oblivious Message Retrieval

## 1.2.2 Our Techniques

**SIMD-Aware Homomorphic Compression.** Our main technical contribution is identifying that the bottleneck of prior works [LT22, LTW24b] is at homomorphic compression and devising a new mechanism based on the SIMD-aware homomorphic compression scheme of [CLPY24]. Prior approaches rely on a *hashing-based* compression scheme (Sec. 3.4) in which each index is *hashed* into *buckets*, making the scheme incompatible with the homomorphic SIMD structure. As a result, before the actual compression, each slot in the detection result must be *unpacked* into separate ciphertexts, using costly subroutines like SlotToCoeff [GHS12a, LW23] and *oblivious expansion* [ACLS18]. The latter step involves  $O(N)$  homomorphic rotations, where  $N$  is the size of the bulletin board, creating a significant bottleneck in the digesting process.

In contrast, our SIMD-aware homomorphic compression mechanism based on [CLPY24] eliminates the need for a complex and expensive unpacking process, enabling a seamless transition between the detection and compression phase (Sec. 4.1). Furthermore, our compression scheme ultimately reduces to a single homomorphic matrix-vector multiplication, a well-studied operation in the HE literature [HS14, HS18, JVC18, CDKS19, LHH<sup>+</sup>21]. As a result, our scheme is arguably simpler to understand and implement, as well as significantly more efficient: it requires only  $O(\sqrt{N})$  homomorphic rotations, compared to the  $O(N)$  rotations required by previous hashing-based approaches. Another advantage of our compression scheme is that, unlike the previous hashing-based approach, it is *deterministic* and does not introduce any additional failure probability, which substantially simplifies the analysis for parameter setting.

A limitation of the SIMD-aware homomorphic compression of [CLPY24] is that it only considers small payloads that fit within a single message *slot* of the underlying HE (e.g., 16-bit). However, in many applications, including OMR, payloads are significantly larger and require multiple slots.

To address this, we introduce *precomputation* and *stacking* techniques to accelerate the compression of multi-slot payloads, achieving a quadratic reduction in the required number of homomorphic rotations, relative to the number of slots, compared to the naive approach of compressing each slot individually.

**Optimizations.** We also suggest several other optimization techniques to further improve the performance of our scheme.

- *Ring-Switching* (Sec. 4.2): We adopt *ring-switching* [BGV12, GHPS12, GHPS13], which reduces the HE parameter after completing the desired homomorphic computations, to achieve further compression of the digest.
- *Hybrid Key-Switching Optimization* (Sec. 4.3): We optimize the parameter for *hybrid key-switching* [GHS12b, HK20, KPZ21] in the HE scheme to reduce the key size.
- *BSGS Optimization for Affine Transform* (Sec. 4.4): We apply *BSGS (Baby-Step-Giant-Step)* optimization to speed up homomorphic matrix-vector multiplications.

## 1.3 Related Work

**Oblivious Message Retrieval.** Besides the two works [LT22, LTW24b] previously mentioned, OMR has been extended further by several studies. Group OMR [LTW24a] extends OMR to group messaging scenarios, offering enhanced efficiency than naive use of OMR, where the detector’s workload scales linearly with the group size. DoS-resistant OMR [LSTW24] considers an extended security notion to safeguard OMR schemes against spamming attacks by malicious senders. We believe that our SophOMR scheme is compatible with techniques developed for the group setting [LTW24a] and DoS-resistance [LSTW24]; however, these extensions are beyond the scope of this work.

**Receiver Privacy.** Several non-HE-based solutions have been proposed that offer receiver privacy along with OMR/OMD-like functionality. Fuzzy Message Detection (FMD) [BLMG21] is a *decoy-based* solution that employs a detection algorithm specifically designed to misidentify impertinent messages at a predetermined false positive rate, effectively hiding the exact set of pertinent messages. However, the decoy-based security notion of FMD is indeed relatively complex and weak, as analyzed in [SPB22]. In contrast to FMD, some works achieve full privacy [MSS<sup>+</sup>22, JLM23, JMK24], but they rely on stronger trust assumptions than OMR, such as trusted execution environment (TEE) or two non-colluding servers.

**Homomorphic Compression.** Several homomorphic compression schemes and its variants have been proposed [CDG<sup>+</sup>21, LT22, FLS23, BPSY24, CLPY24]. However, many of these schemes [CDG<sup>+</sup>21, LT22, FLS23] are not particularly efficient when implemented with HE schemes, as they are not SIMD-friendly due to their *hashing-based* approach. (See Sec. 3.4.2.) Although Cheon-Lee-Park-Yeo [CLPY24] is the first to investigate SIMD-friendliness, some earlier works feature structures that can be implemented in a SIMD-friendly manner. However, these approaches have their own limitations: The FFT-based scheme by Fleischhacker-Larsen-Simkin [FLS23] has a costly decompression algorithm that requires solving small instances of discrete logarithms, and the scheme by Bienstock-Patel-Seo-Yeo [BPSY24] assumes that the entity performing the decompression already knows the indices of interest, which is not applicable in the context of OMR. Cheon et al. [CLPY24] proposed a SIMD-aware homomorphic compression scheme by reinterpreting and extending PS-COIE of Choi et. al. [CDG<sup>+</sup>21], which only considered compressing binary vector and disregarded its SIMD-friendly structure. This new scheme achieves both an asymptotically optimal compression rate and asymptotically good decompression complexity.

**Comparison to PIR.** Private Information Retrieval (PIR) [CGKS95, KO97] is a cryptographic primitive that allows clients to retrieve an item from a server hosting a database without revealing which specific item is being accessed. Although PIR and OMR might appear similar, they exhibit fundamentally different characteristics. In PIR, each query is made by an index that must remain hidden from the server, whereas in OMR, queries are executed using the receiver’s public key, which does not necessarily need to be hidden. Furthermore, a typical PIR scheme retrieves a single entry, but an OMR scheme retrieves multiple pertinent payloads. Most importantly, an OMR scheme incorporates a mechanism that allows anyone to tag their messages with a *signal*, which can only be detected by the intended receiver.

## 2 Preliminaries

### 2.1 Notations

We use  $[n]$  to denote the set of integers  $\{0, \dots, n-1\}$ , and  $\log$  refers to the base-2 logarithm. Throughout the paper, we follow the column vector notation. For a vector  $\mathbf{v}$ , we use  $\mathbf{v}[i]$  to denote its  $i$ -th element, and for a matrix  $\mathbf{M}$ , we use  $\mathbf{M}[i][j]$  to denote its  $(i, j)$ -th element. The Hadamard product (a.k.a. element-wise product) is denoted by  $\odot$ , and  $\langle \cdot, \cdot \rangle$  denotes the inner product. For  $i < n$ , we use  $\text{Rot}^i$  to denote the left rotation of a vector by  $i$  positions, i.e.,  $\text{Rot}^i(v_0, \dots, v_{n-1}) = (v_i, \dots, v_{n-1}, v_0, \dots, v_{i-1})$ . The vector  $\mathbf{1}$  refers to a vector of ones  $(1, \dots, 1)$ . We use  $\lambda$  for the security parameter and  $\text{negl}(\lambda)$  to denote a negligible function with respect to  $\lambda$ . The notation  $\text{ct}(m)$  denotes a ciphertext encrypting a message  $m$ , and  $\text{ct}(\mathbf{m})$  denotes possibly multiple ciphertexts encrypting the entries of a vector  $\mathbf{m}$ .

### 2.2 Homomorphic Encryption

A homomorphic encryption (HE) scheme allows us to compute an encryption of  $f(m)$  when only given a circuit  $f$  and a ciphertext encrypting  $m$  without decryption. Let  $\mathcal{M}$  be the message space and  $\mathcal{C}$  be the ciphertext space. An HE scheme is a tuple of probabilistic polynomial time algorithms (KeyGen, Enc, Dec, Eval) described below:

- $\text{KeyGen}(1^\lambda)$ : The key generation algorithm takes a security parameter  $\lambda$  and returns a secret key  $\text{sk}$  and public key  $\text{pk}$ .
- $\text{Enc}_{\text{pk}}(m)$ : The encryption algorithm takes the public key  $\text{pk}$  and a message  $m \in \mathcal{M}$  as inputs. Then, it returns a ciphertext  $\text{ct}(m) \in \mathcal{C}$  encrypting  $m$ .
- $\text{Dec}_{\text{sk}}(\text{ct}(m))$ : The decryption algorithm takes the secret key  $\text{sk}$  and a ciphertext  $\text{ct}(m) \in \mathcal{C}$  encrypting a message  $m \in \mathcal{M}$ . Then, it returns the message  $m$ .
- $\text{Eval}_{\text{pk}}(f, \text{ct}(\mathbf{m}))$ : The evaluation algorithm takes the public key  $\text{pk}$ , a circuit  $f$ , and ciphertext  $\text{ct}(\mathbf{m}) \in \mathcal{C}^*$  as inputs. Then, it returns ciphertext  $\text{ct}(f(\mathbf{m})) \in \mathcal{C}^*$  that encrypts  $f(\mathbf{m}) \in \mathcal{M}^*$  as message.

**Definition 1** (Security of HE). *Let  $(\text{sk}, \text{pk}) \leftarrow \text{KeyGen}(1^\lambda)$ .*

- **Correctness:** *For any circuit  $f$  (in some class  $\mathcal{F}$ ) and input  $\mathbf{m} \in \mathcal{M}^*$ , the following holds with overwhelming probability for  $\text{ct}(\mathbf{m}) \leftarrow \text{Enc}_{\text{pk}}(\mathbf{m})$ :*

$$\text{Dec}_{\text{sk}}(\text{Eval}_{\text{pk}}(f, \text{ct}(\mathbf{m}))) = f(\mathbf{m})$$

- **Semantic Security:** *For any probabilistic polynomial time adversary  $\mathcal{A}$  and  $m_1, m_2 \in \mathcal{M}$ , the following holds for  $\text{ct}_i \leftarrow \text{Enc}_{\text{pk}}(m_i)$ :*

$$|\Pr[\mathcal{A}(\text{pk}, \text{ct}_1) = 1] - \Pr[\mathcal{A}(\text{pk}, \text{ct}_2) = 1]| \leq \text{negl}(\lambda)$$

### 2.2.1 BFV Scheme

As in the previous works [LT22, LTW24b], we employ the BFV scheme [Bra12, FV12] in our SophOMR construction.<sup>7</sup> The BFV scheme supports homomorphic SIMD arithmetic in  $\mathbb{Z}_p$  for a prime  $p$ , which is critical to our SophOMR construction, especially in our homomorphic compression scheme.

More precisely, we consider the BFV scheme with the ring dimension  $n$  that is a power of 2, along with the plaintext modulus  $p$  satisfying  $p = 1 \pmod{2n}$ . In this case, the plaintext space is given by  $\mathcal{R}_p := \mathbb{Z}_p[X]/(X^n + 1) \cong \mathbb{Z}_p^{2 \times (n/2)}$ , while the ciphertext space is  $\mathcal{R}_Q^2$ , where  $\mathcal{R}_Q := \mathbb{Z}_Q[X]/(X^n + 1)$  for a ciphertext modulus  $Q$ . For a plaintext  $\mathbf{m} \in \mathbb{Z}_p^{2 \times (n/2)}$ , we denote its rows and entries as follows. We often refer to each plaintext entry as a *slot*.

$$\mathbf{m} = \begin{bmatrix} \mathbf{m}_1^T \\ \mathbf{m}_2^T \end{bmatrix} = \begin{bmatrix} m_{1,1} & \cdots & m_{1,n/2} \\ m_{2,1} & \cdots & m_{2,n/2} \end{bmatrix} \quad (1)$$

The BFV scheme supports (1) SIMD arithmetic in  $\mathbb{Z}_p$  and (2) homomorphic rotations (horizontal and vertical). That is, we can homomorphically add and multiply a ciphertext encrypting  $\mathbf{v} \in \mathbb{Z}_p^{2 \times (n/2)}$  with a ciphertext (or plaintext) of  $\mathbf{w} \in \mathbb{Z}_p^{2 \times (n/2)}$  to obtain a ciphertext of  $\mathbf{v} + \mathbf{w}$  and  $\mathbf{v} \odot \mathbf{w}$ , respectively. For homomorphic rotations, to be precise, let  $\text{Rot}^r(\cdot)$  denote the rotation operation which maps  $[v_1, \dots, v_{n/2}] \mapsto [v_{1+r}, \dots, v_{n/2}, v_1, \dots, v_r]$ . Then, the BFV scheme supports homomorphic rotations computing  $\text{Rot}_{\text{row}}$  and  $\text{Rot}_{\text{col}}$ , defined as follows.

$$\begin{aligned} \text{Rot}_{\text{row}}^r \left( \begin{bmatrix} \mathbf{m}_1^T \\ \mathbf{m}_2^T \end{bmatrix} \right) &= \begin{bmatrix} \text{Rot}^r(\mathbf{m}_1^T) \\ \text{Rot}^r(\mathbf{m}_2^T) \end{bmatrix} \\ \text{Rot}_{\text{col}}^1 \left( \begin{bmatrix} \mathbf{m}_1^T \\ \mathbf{m}_2^T \end{bmatrix} \right) &= \begin{bmatrix} \mathbf{m}_2^T \\ \mathbf{m}_1^T \end{bmatrix} \end{aligned}$$

In this work, we will often abuse notation and identify a BFV plaintext  $[\mathbf{m}_1 | \mathbf{m}_2]^T \in \mathbb{Z}_p^{2 \times (n/2)}$  with the vector  $[\mathbf{m}_1^T | \mathbf{m}_2^T]^T \in \mathbb{Z}_p^n$ .

### 2.2.2 Sparse Packing

When working with a vector  $\mathbf{v}$  of length  $k$  that is smaller than the HE ring dimension  $n$ , we often employ *Sparse Packing* (a.k.a. Repetitive Slot Packing), where multiple copies of  $\mathbf{v}$  are packed into a plaintext. More precisely, when  $n = 0 \pmod{k}$ , we pack  $\mathbf{v}$  ( $n/k$ ) times, resulting in  $[\mathbf{v}^T | \dots | \mathbf{v}^T]^T$ . The rationale for using sparse packing is that it preserves the rotation structure. For example,  $\text{Rot}_{\text{row}}^k([\mathbf{v}^T | \dots | \mathbf{v}^T]^T) = [\mathbf{v}^T | \dots | \mathbf{v}^T]^T$ , which would not hold if we did not use sparse packing, such as when padding the remaining slots naively with zeros.

<sup>7</sup>We note that SophOMR is also compatible with the BGV scheme [BGV12]. Although one might consider using the TFHE scheme [CGG20], which supports homomorphic operations without requiring SIMD structure, it incurs significant communication costs even after compression due to the high plaintext-ciphertext expansion ratio.

Throughout this paper, whenever  $n = 0 \pmod{\dim \mathbf{v}}$ ,  $\text{ct}(\mathbf{v})$  should be interpreted as a sparsely packed ciphertext.

## 2.3 Homomorphic Matrix Multiplication

At the heart of our SophOMR scheme is homomorphic matrix-vector multiplication (MatMul) [HS14, HS18, JVC18, CDKS19, LHH<sup>+</sup>21]. The goal of MatMul is to multiply a plaintext matrix  $\mathbf{M}$  with an encryption of vector  $\mathbf{v}$  and obtain an encryption of  $\mathbf{M}\mathbf{v}$ . Here, we review an efficient MatMul algorithm that we leverage in our SophOMR scheme.

### 2.3.1 Diagonal Packing

The idea of *diagonal packing*, first introduced in [HS14], is fundamental to MatMul algorithms. Given a matrix  $\mathbf{M} \in \mathbb{Z}_p^{m \times k}$  and a vector  $\mathbf{v} \in \mathbb{Z}_p^k$ , let  $\text{diag}_i(\mathbf{M})$  denote the  $i$ -th diagonal packing of  $\mathbf{M}$  (w.r.t. dimension  $n$ ), which is defined as follows (for  $0 \leq j < n$ ).

$$\text{diag}_i(\mathbf{M})[j] = \mathbf{M}[j \bmod m][[(i+j) \bmod k]]$$

The rationale behind diagonal packing is that it enables a concise representation of a matrix-vector multiplication with native HE operations, specifically Hadamard products and rotations, as follows (for  $n = m = k$ ).

$$\mathbf{M}\mathbf{v} = \sum_{i=0}^{k-1} \text{diag}_i(\mathbf{M}) \odot \text{Rot}^i(\mathbf{v}) \quad (2)$$

### 2.3.2 BSGS-Style MatMul

Halevi and Shoup [HS18] introduced the baby-step-giant-step (BSGS) style optimization to reduce the number of homomorphic rotations required in a MatMul algorithm. The core idea is captured by the following equalities, where  $k = \tilde{g} \cdot \tilde{b}$  and  $\mathbf{m}_{g\tilde{b}+b}$  denotes  $\text{Rot}^{-g\tilde{b}}(\text{diag}_{g\tilde{b}+b}(\mathbf{M}))$ . Notice that, by reusing the  $\text{Rot}^b(\mathbf{v})$  terms, we can reduce the number of rotations to  $\tilde{g} + \tilde{b}$ , rather than  $k = \tilde{g} \cdot \tilde{b}$  as in Eq. 2.

$$\begin{aligned} \mathbf{M}\mathbf{v} &= \sum_{i=0}^{k-1} \text{diag}_i(\mathbf{M}) \odot \text{Rot}^i(\mathbf{v}) \\ &= \sum_{g=0}^{\tilde{g}-1} \sum_{b=0}^{\tilde{b}-1} \text{diag}_{g\tilde{b}+b}(\mathbf{M}) \odot \text{Rot}^{g\tilde{b}+b}(\mathbf{v}) \\ &= \sum_{g=0}^{\tilde{g}-1} \text{Rot}^{g\tilde{b}} \left( \sum_{b=0}^{\tilde{b}-1} \mathbf{m}_{g\tilde{b}+b} \odot \text{Rot}^b(\mathbf{v}) \right) \end{aligned}$$

Among several follow-ups on [HS18], we adopt the approach of PEGASUS [LHH<sup>+</sup>21], which exhibits the best performance in the case of non-square matrices. The original MatMul algorithm of PEGASUS is designed for the CKKS scheme [CKKS17] and its rotation structure. In this

work, we use its straightforward adaptation to our BFV setting [Bra12, FV12]. Additionally, we extend the approach to handle cases where the matrix dimension exceeds the underlying HE ring dimension, requiring the encrypted vectors to span multiple ciphertexts.

To be precise, let  $n$  be the ring dimension of the underlying BFV scheme. Let us consider the MatMul of  $k \times N$  matrix  $\mathbf{M}$  with a length- $N$  vector  $\mathbf{v}$ . Here, we focus on the *wide* case, where  $N = s \cdot n$  for some positive integer  $s$  and  $k < n$ .<sup>8</sup> For simplicity, we assume  $k$  is a power of two. Alg. 1 describes the BFV-adapted PEGASUS MatMul algorithm for the specified setting. We use the Horner-style [HS18] variant, which optimizes the number of required rotation keys.<sup>9</sup> The costs for Alg. 1 are summarized in Tab. 1. Note that the number of rotations is minimized when  $\tilde{g} \approx \sqrt{s \cdot k}$ . For more detailed discussions, refer to PEGASUS [LHH<sup>+</sup>21] and references therein, such as [HS18].

---

**Algorithm 1** MatMul (Wide)

---

**Input:**  $\mathbf{M} \in \mathbb{Z}_p^{k \times N}$ ,  $(\text{ct}(\mathbf{v}_i))_{i=0}^{s-1} \triangleright k < n, N = s \cdot n, \mathbf{v}_i \in \mathbb{Z}_p^n$ ,  
 $k$ : power-of-two

**Step 0:** Parameter Setting

$$[\mathbf{M}_0 | \dots | \mathbf{M}_{s-1}] \leftarrow \mathbf{M} \quad \text{and} \quad k = \tilde{g} \cdot \tilde{b}$$

**Step 1:** Plaintext Packing

**for**  $(i, j) \in [s] \times [k]$  **do**

$$\mathbf{m}_{i,j} \leftarrow \text{Rot}_{\text{row}}^{-g\tilde{b}}(\text{diag}_j(\mathbf{M}_i)) \quad \text{for } g = \lfloor j/\tilde{b} \rfloor$$

**Step 2:** Baby-Step-Giant-Step (BSGS)

**for**  $i = 0, \dots, s-1$  **do**  $\triangleright$  Baby-Step

$$\text{ct}_{i,0} \leftarrow \text{ct}(\mathbf{v}_i)$$

**for**  $b = 1, \dots, \tilde{b}-1$  **do**

$$\text{ct}_{i,b} \leftarrow \text{Rot}_{\text{row}}^1(\text{ct}_{i,b-1})$$

**for**  $g = \tilde{g}-1, \dots, 0$  **do**  $\triangleright$  Giant-Step

$$\text{ct}_{\text{sum}} \leftarrow \sum_{i=0}^{s-1} \sum_{b=0}^{\tilde{b}-1} (\mathbf{m}_{i,g\tilde{b}+b} \odot \text{ct}_{i,b})$$

**if**  $g = \tilde{g}-1$  **then**

$$\text{ct}_{\text{out}} \leftarrow \text{ct}_{\text{sum}}$$

**else**

$$\text{ct}_{\text{out}} \leftarrow \text{Rot}_{\text{row}}^{\tilde{b}}(\text{ct}_{\text{out}}) + \text{ct}_{\text{sum}}$$

**Step 3:** Block Summation

**for**  $j = 0, \dots, \log(n/k) - 2$  **do**

$$\text{ct}_{\text{out}} \leftarrow \text{ct}_{\text{out}} + \text{Rot}_{\text{row}}^{k \cdot 2^j}(\text{ct}_{\text{out}})$$

$$\text{ct}_{\text{out}} \leftarrow \text{ct}_{\text{out}} + \text{Rot}_{\text{col}}^1(\text{ct}_{\text{out}})$$

**return**  $\text{ct}_{\text{out}}$

---

<sup>8</sup>The *tall* case is described in Alg. 8 of Sec. 4.4.

<sup>9</sup>While Horner-style MatMul algorithms benefit from the significantly reduced number of rotation keys, they are incompatible with so-called *hoisting* optimization [HS18]. This trade-off is discussed in Sec. 4.1.

#(Ptxt-Mult)	#(Rotation)	#(Rot. Key)
$s \cdot k$	$< s \cdot \tilde{b} + \tilde{g} + \log(n/k)$	$\leq 2 + \log(n/k)$

Table 1: Cost Analysis for Alg. 1

## 2.4 Homomorphic Compression

The goal of homomorphic compression [CDG<sup>+</sup>21, LT22, FLS23, BPSY24, CLPY24] is to homomorphically compress an encryption of a *sparse* vector, which has only a few non-zero components, into an encryption of a shorter vector, which can later be efficiently decompressed after decryption. In this section, we provide a brief overview of the homomorphic compression scheme from [CLPY24], which offers high efficiency due to its compatibility with homomorphic SIMD operations. For a more in-depth discussion, please refer to [CLPY24].

### 2.4.1 Index Compression

We first consider the task of compressing a sparse *binary* vector (a.k.a. index compression) as a stepping stone for compressing an arbitrary sparse vector (a.k.a. payload compression). The index compression of [CLPY24] leverages so-called *Newton's identity*, which provides conversions between *power sum polynomials* and *elementary symmetric polynomials*. Although the idea of using Newton's identity for index compression has been previously explored in the literature [CDG<sup>+</sup>21], the novelty of [CLPY24] lies in (1) recognizing that the power sum encoding can be interpreted as a matrix-vector multiplication and (2) implementing it with efficient SIMD-aware homomorphic algorithms for MatMul (e.g., Alg. 1).

More precisely, for  $p > N$ , consider a sparse binary vector  $\mathbf{v} \in \mathbb{Z}_p^N$  with at most  $k$  number of non-zero entries and define a Vandermonde-like  $k \times N$  matrix  $\mathbf{C}$  defined as follows. The homomorphic compression scheme can then be concisely described as in Alg. 2, with decompression outlined in Alg. 3.

$$\mathbf{C} := \begin{bmatrix} 1 & 2 & \dots & N \\ 1^2 & 2^2 & \dots & N^2 \\ \vdots & \vdots & \ddots & \vdots \\ 1^k & 2^k & \dots & N^k \end{bmatrix} \pmod{p} \quad (3)$$

---

**Algorithm 2** Complx

---

**Input:**  $\text{ct}(\mathbf{v})$

$\triangleright$  ciphertexts of  $\mathbf{v} \in \mathbb{Z}_p^N$

$$\text{ct}(\mathbf{w}) \leftarrow \text{MatMul}(\mathbf{C}, \text{ct}(\mathbf{v}))$$

$\triangleright$  see Sec. 2.3

**return**  $\text{ct}(\mathbf{w})$

---

---

**Algorithm 3** Decompldx

---

**Input:**  $\mathbf{w} = (w_1, \dots, w_k) \in \mathbb{Z}_p^k$   $\triangleright$  decryption of  $\text{ct}(\mathbf{w})$

**Step 1:** Build a polynomial from  $\mathbf{w}$  using Newton's identity.

$a_0 \leftarrow 1$   
**for**  $j = 1, \dots, k$  **do**  
     $a_j \leftarrow j^{-1} \cdot \sum_{i=1}^j (-1)^{i-1} a_{j-i} \cdot w_i$   
 $f_{\mathbf{v}}^k(X) \leftarrow \sum_{i=0}^k (-1)^i a_i X^{k-i}$

**Step 2:** Reconstruct non-zero indices from the polynomial.

**while**  $X$  divides  $f_{\mathbf{v}}^k(X)$  **do**  
     $f_{\mathbf{v}}^k(X) \leftarrow f_{\mathbf{v}}^k(X)/X$   
 $I_{\mathbf{v}} \leftarrow$  Find zeros of  $f_{\mathbf{v}}^k(X)$   $\triangleright$  e.g. by Cantor-Zassenhaus

**return**  $I_{\mathbf{v}}$

---

### 2.4.2 Payload Compression

With the index compression at hand, we can now compress an arbitrary sparse vector. Consider a sparse vector  $\mathbf{d} \in \mathbb{Z}_p^N$  with at most  $k$  number of non-zero entries. To compress  $\text{ct}(\mathbf{d})$ , we first obtain an encryption of the binary vector  $\mathbf{v}$ , whose indices of non-zero entries match those in  $\mathbf{d}$ . For example, Fermat's little theorem can be employed for this purpose (as outlined in Alg. 4). However, in certain scenarios (e.g., private database query [CLPY24]), it may not be necessary to derive  $\text{ct}(\mathbf{v})$  from  $\text{ct}(\mathbf{d})$ , as  $\text{ct}(\mathbf{v})$  may be already available from other source. In fact,  $\text{ct}(\mathbf{d})$  is often computed *from*  $\text{ct}(\mathbf{v})$ . Looking ahead, this is the case also for OMR.

Then, we can compress  $\text{ct}(\mathbf{v})$  via Alg. 2 and use  $\mathbf{v}$  as a *hint* in the decompression phase. Note that any  $k$  distinct columns of the matrix  $\mathbf{C}$  are linearly independent. Therefore, we can also compress  $\text{ct}(\mathbf{d})$  by performing a MatMul with  $\mathbf{C}$  and decompress it by solving a linear equation for a submatrix of  $\mathbf{C}$ . The homomorphic compression scheme can then be concisely described as in Alg. 4, with decompression outlined in Alg. 5.

---

**Algorithm 4** Comp

---

**Input:**  $\text{ct}(\mathbf{d})$   $\triangleright$  ciphertexts of  $\mathbf{d} \in \mathbb{Z}_p^N$

$\text{ct}(\mathbf{v}) \leftarrow \text{Eval}(\text{Power}_{p-1}, \text{ct}(\mathbf{d}))$ <sup>10</sup>  $\triangleright$  Fermat's little theorem  
 $\text{ct}(\mathbf{w}) \leftarrow \text{Compldx}(\text{ct}(\mathbf{v}))$   $\triangleright$  see Alg. 2  
 $\text{ct}(\mathbf{e}) \leftarrow \text{Eval}(\mathbf{C}, \text{ct}(\mathbf{d}))$

**return**  $\text{ct}(\mathbf{e}), \text{ct}(\mathbf{w})$

---

## 2.5 OMR and OMD

In this section, we review the formal definitions of Oblivious Message Retrieval (OMR) and Oblivious Message Detection (OMD) [LT22]. See also Sec. 1.2.1 and Fig. 1.

<sup>10</sup>Power $_{\ell}$  denotes an arithmetic circuit over  $\mathbb{Z}_p$  such that Power $_{\ell}(\mathbf{d}) = (d_1^{\ell}, \dots, d_N^{\ell})$ , where  $\mathbf{d} = (d_1, \dots, d_N) \in \mathbb{Z}_p^N$ .

---

**Algorithm 5** Decomp

---

**Input:**  $\mathbf{e}, \mathbf{w} \in \mathbb{Z}_p^k$   $\triangleright$  decryption of  $\text{ct}(\mathbf{e}), \text{ct}(\mathbf{w})$

$I_{\mathbf{v}} \leftarrow \text{Decompldx}(\mathbf{w})$   $\triangleright$  see Alg. 3  
 $\hat{\mathbf{C}} \leftarrow [\mathbf{C}_{i_1} | \dots | \mathbf{C}_{i_{\ell}}]$   $\triangleright k \times \ell$  matrix whose  $j$ -th column is  $i_j$ -th column of  $\mathbf{C}$ , where  $I_{\mathbf{v}} = \{i_1, \dots, i_{\ell}\}$

$\hat{\mathbf{d}} = (\hat{d}_1, \dots, \hat{d}_{\ell}) \leftarrow$  Solve  $\hat{\mathbf{C}}\mathbf{x} = \mathbf{e}$  for  $\mathbf{x}$

**return**  $\hat{\mathbf{d}}$

---

### 2.5.1 Oblivious Message Retrieval (OMR)

**Definition 2** (Oblivious Message Retrieval). *An Oblivious Message Retrieval (OMR) scheme (w.r.t. payload space  $\mathcal{P}$ ) is a tuple of probabilistic algorithms (ParamGen, KeyGen, Signal, Digest, Decode) described below:*

- $\text{pp} \leftarrow \text{ParamGen}(1^{\lambda})$ : *The parameter generation algorithm takes a security parameter  $\lambda$  and returns a public parameter  $\text{pp}$ . The following algorithms implicitly take  $\text{pp}$  as an input.*
- $(\text{sk}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}}) \leftarrow \text{KeyGen}(\text{pp})$ : *The key generation algorithm outputs a secret key  $\text{sk}$ , a signal key  $\text{pk}_{\text{sig}}$ , and a detection key  $\text{pk}_{\text{det}}$ .*
- $\text{sig} \leftarrow \text{Signal}(\text{pk}_{\text{sig}})$ : *The signaling algorithm takes a signal key  $\text{pk}_{\text{sig}}$  and outputs a signal  $\text{sig}$ .*

*A bulletin board  $\text{BB} = ((x_1, \text{sig}_1), \dots, (x_N, \text{sig}_N))$  is formed by several parties uploading a payload  $x_i \in \mathcal{P}$  tagged with a signal  $\text{sig}_i$  with respect to the signal key of the intended receiver.*

- $\text{D} \leftarrow \text{Digest}(\text{BB}, \text{pk}_{\text{det}}, k)$ : *The digesting algorithm takes a bulletin board  $\text{BB} = ((x_1, \text{sig}_1), \dots, (x_N, \text{sig}_N))$  for some size  $N$ , a detection key  $\text{pk}_{\text{det}}$ , and an upper bound  $k$  on the number of pertinent payloads addressed to the receiver. Then, it returns a digest  $\text{D}$ .*
- $\text{M} \leftarrow \text{Decode}(\text{D}, \text{sk})$ : *The decoding algorithm takes a digest  $\text{D}$  and a secret key  $\text{sk}$ . Then, it returns a list of decoded payloads  $\text{M} \in \mathcal{P}^*$ .*

**Definition 3** (Security of OMR). *Let  $\text{pp} \leftarrow \text{ParamGen}(1^{\lambda})$ .*

- **Correctness:** *Let  $(\text{sk}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}}) \leftarrow \text{KeyGen}(\text{pp})$ . Let  $\text{BB}$  be a bulletin board of length  $N$  and  $\tilde{\text{M}}$  be the list of pertinent payloads in  $\text{BB}$  addressed to  $\text{pk}_{\text{sig}}$ . For  $k \geq |\tilde{\text{M}}|$ , the following holds.*

$$\Pr \left[ \text{M} \neq \tilde{\text{M}} \mid \begin{array}{l} \text{D} \leftarrow \text{Digest}(\text{BB}, \text{pk}_{\text{det}}, k) \\ \text{M} \leftarrow \text{Decode}(\text{D}, \text{sk}) \end{array} \right] \leq \text{negl}(\lambda)$$

- **Receiver-Privacy:** *Signals generated from different signal keys should be computationally indistinguishable to anyone who does not have access to the corresponding secret keys.*

To be precise, let  $(\text{sk}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}})$  and  $(\text{sk}', \text{pk}'_{\text{sig}}, \text{pk}'_{\text{det}})$  be two set of keys generated from  $\text{KeyGen}(\text{pp})$ . Let  $\text{sig} \leftarrow \text{Signal}(\text{pk}_{\text{sig}})$  and  $\text{sig}' \leftarrow \text{Signal}(\text{pk}'_{\text{sig}})$ . For any probabilistic polynomial time adversary  $\mathcal{A}$ , the following holds for  $\text{aux} = (\text{pp}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}}, \text{pk}'_{\text{sig}}, \text{pk}'_{\text{det}})$ .

$$|\Pr[\mathcal{A}(\text{aux}, \text{sig}) = 1] - \Pr[\mathcal{A}(\text{aux}, \text{sig}') = 1]| \leq \text{negl}(\lambda)$$

### 2.5.2 Oblivious Message Detection (OMD)

*Oblivious Message Detection (OMD)* is a variant of OMR where the receiver only retrieves pertinent *indices* (instead of pertinent *payloads* as in OMR). Note that by combining an OMD scheme with a private information retrieval (PIR) scheme [CGKS95, KO97] one can obtain an OMR-like functionality: first, retrieve pertinent indices using OMD, and then use PIR to obtain the corresponding payloads. While this approach introduces an additional round of communication, it provides a generic construction of OMR. Below are formal definitions for OMD.

**Definition 4** (Oblivious Message Detection). *An Oblivious Message Detection (OMD) scheme is a tuple of probabilistic algorithms  $(\text{ParamGen}, \text{KeyGen}, \text{Signal}, \text{Digest}, \text{Decode})$ . Descriptions for  $\text{ParamGen}$ ,  $\text{KeyGen}$ ,  $\text{Signal}$ , and  $\text{Digest}$  are identical to those of OMR (Def. 2).*

- $\text{M} \leftarrow \text{Decode}(\text{D}, \text{sk})$ : *The decoding algorithm takes a digest  $\text{D}$  and a secret key  $\text{sk}$ . Then, it returns a list of pertinent indices  $\text{M} \in \mathbb{N}^*$ .*

**Definition 5** (Security of OMD). *Let  $\text{pp} \leftarrow \text{ParamGen}(1^\lambda)$ . The definition of receiver-privacy is identical to that of OMR (Def. 3).*

- **Correctness:** *Let  $(\text{sk}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}}) \leftarrow \text{KeyGen}(\text{pp})$ . Let  $\text{BB}$  be a bulletin board of length  $N$  and  $\tilde{\text{M}}$  be the list of pertinent indices addressed to  $\text{pk}_{\text{sig}}$  in  $\text{BB}$ . For  $k \leq |\tilde{\text{M}}|$ , the following holds.*

$$\Pr \left[ \text{M} \neq \tilde{\text{M}} \mid \begin{array}{l} \text{D} \leftarrow \text{Digest}(\text{BB}, \text{pk}_{\text{det}}, k) \\ \text{M} \leftarrow \text{Decode}(\text{D}, \text{sk}) \end{array} \right] \leq \text{negl}(\lambda)$$

## 3 Review: PerfOMR [LTW24b]

In this section, we provide a concise overview of the PerfOMR scheme [LTW24b]. First, in Sec. 3.1, we recall a building block for OMR, which we call *private signaling*. Then, we review each component of PerfOMR:  $\text{ParamGen}$ ,  $\text{KeyGen}$ ,  $\text{Signal}$  (Sec. 3.2), and  $\text{Digest}$  (Sec. 3.3 and 3.4). Specifically, the  $\text{Digest}$  algorithm is presented in a modular manner, split into a detection phase (Sec. 3.3) and a compression phase (Sec. 3.4).

## 3.1 RLWE-based Private Signaling Scheme

A Private Signaling (PS) scheme is a core building block of OMR, enabling a sender to generate a *signal* that only the intended receiver can detect. We begin by presenting a formal definition of private signaling (Sec. 3.1.1), followed by a construction based on the Ring Learning with Error (RLWE) assumption [Reg05, LPR10] (Sec. 3.1.2).

### 3.1.1 Private Signaling

**Definition 6** (Private Signaling<sup>11</sup>). *A Private Signaling (PS) scheme is a tuple of probabilistic polynomial time algorithms  $(\text{ParamGen}, \text{KeyGen}, \text{Signal}, \text{Detect})$  described below:*

- $\text{pp} \leftarrow \text{ParamGen}(1^\lambda)$ : *The parameter generation algorithm takes a security parameter  $\lambda$  and returns a public parameter  $\text{pp}$ . The following algorithms implicitly take  $\text{pp}$  as an input.*
- $(\text{sk}, \text{pk}) \leftarrow \text{KeyGen}(\text{pp})$ : *The key generation algorithm returns a secret key  $\text{sk}$  and a public key  $\text{pk}$ .*
- $\text{sig} \leftarrow \text{Signal}(\text{pk})$ : *The signaling algorithm takes a public key  $\text{pk}$  and returns a signal  $\text{sig}$ .*
- $0/1 \leftarrow \text{Detect}_{\text{sk}}(\text{sig})$ : *The detection algorithm takes a secret key  $\text{sk}$  and a signal  $\text{sig}$ . Then, it returns 1 if  $\text{sig}$  was generated from a public key corresponding to  $\text{sk}$ ; otherwise, it returns 0.*

**Definition 7** (Security of Private Signaling). *Let  $\text{pp} \leftarrow \text{ParamGen}(1^\lambda)$ .*

- **Completeness:** *The following holds.*

$$\Pr \left[ \begin{array}{l} (\text{sk}, \text{pk}) \leftarrow \text{KeyGen}(\text{pp}); \\ \text{sig} \leftarrow \text{Signal}(\text{pk}); \\ \text{Detect}_{\text{sk}}(\text{sig}) = 1 \end{array} \right] \geq 1 - \text{negl}(\lambda)$$

- **Soundness:** *The following holds.*

$$\Pr \left[ \begin{array}{l} (\text{sk}, \text{pk}) \leftarrow \text{KeyGen}(\text{pp}); \\ (\text{sk}', \text{pk}') \leftarrow \text{KeyGen}(\text{pp}); \\ \text{sig} \leftarrow \text{Signal}(\text{pk}); \\ \text{Detect}_{\text{sk}'}(\text{sig}) = 1 \end{array} \right] \leq \text{negl}(\lambda)$$

- **Receiver-Privacy:** *Signals generated from different public keys should be computationally indistinguishable to anyone who does not have access to the corresponding secret keys. To be precise, let  $(\text{sk}, \text{pk})$  and  $(\text{sk}', \text{pk}')$  be two set of keys generated from  $\text{KeyGen}(\text{pp})$ . Let  $\text{sig} \leftarrow \text{Signal}(\text{pk})$  and  $\text{sig}' \leftarrow \text{Signal}(\text{pk}')$ . For any probabilistic polynomial time adversary  $\mathcal{A}$ , the following holds for  $\text{aux} = (\text{pp}, \text{pk}, \text{pk}')$ .*

$$|\Pr[\mathcal{A}(\text{aux}, \text{sig}) = 1] - \Pr[\mathcal{A}(\text{aux}, \text{sig}') = 1]| \leq \text{negl}(\lambda)$$

<sup>11</sup>Our definition differs from that of [MSS<sup>+</sup>22] and aims to capture the framework they refer to as *naive approach* based on *key-private* PKE.



### 3.1.2 Construction

We now present an overview of the private signaling scheme from PerfOMR [LTW24b] whose security is based on the hardness of the *Ring Learning with Errors (RLWE)* problem [Reg05, LPR10]. For a more detailed discussion of the scheme, refer to [LTW24b].

**Remark 1** (On Choice of LWE-based Scheme). *Although we have a generic conversion from any PKE scheme with reasonable properties, private signaling schemes in this line of works [LT22, LTW24b], are based on (Ring) LWE assumption. The primary reason is that the decryption circuit of LWE-based PKE is simpler and better suited to the computational model of HE than other cryptographic assumptions (e.g., group-based assumptions). Specifically, a large part of the decryption process is just a simple linear algebra over a relatively small modulus.*

We first introduce a notation to ease the presentation. We denote and define the  $i$ -th negacyclic vector of  $\mathbf{a} = (a_0, \dots, a_{n'-1}) \in \mathbb{Z}_{q'}^{n'}$  as follows, for  $0 \leq i < n'$ .

$$\text{Neg}^i(\mathbf{a}) = (a_i, \dots, a_0, -a_{n'-1}, \dots, -a_{i+1})$$

Let  $a(X) = \sum_{i=0}^{n'-1} a_i X^i$  for  $\mathbf{a}$  and similarly for  $b(X)$  and  $\mathbf{b} = (b_0, \dots, b_{n'-1})$ . Notice that  $\langle \text{Neg}^i(\mathbf{a}), \mathbf{b} \rangle$  equals to the  $i$ -th coefficient of  $a(X) \cdot b(X)$  over  $\mathcal{R}_{q'} = \mathbb{Z}_{q'}[X]/(X^{n'} + 1)$ . With this in mind, the following construction can be seen, at a high level, as an RLWE-based PKE scheme specifically designed for encrypting zeroes. It can be viewed as a *ring-optimized* variant of the PVW scheme [PVW08], or as a *truncated* version of the LPR scheme [LPR13].

**Construction 1.** *The private signaling scheme PS = (PS.ParamGen, PS.KeyGen, PS.Signal, PS.Detect) is defined as follows.*

- $\text{pp} \leftarrow \text{PS.ParamGen}(1^\lambda)$ : *The parameter generation algorithm returns  $\text{pp} = (n', q', \ell, r, \chi_s, \chi_e)$ , consisting of polynomial ring dimension  $n'$ , modulus  $q'$ , repetition parameter  $\ell$ , range parameter  $r$ , secret key distribution  $\chi_s$  over  $\mathbb{Z}_{q'}^{n'}$ , and error distribution  $\chi_e$  over  $\mathbb{Z}_{q'}^{\ell}$ .*
- $(\text{sk}, \text{pk}) \leftarrow \text{PS.KeyGen}(\text{pp})$ : *The key generation algorithm first samples a secret key  $\text{sk} \leftarrow \chi_s$ . Then, it samples  $\boldsymbol{\alpha} \leftarrow \mathbb{Z}_{q'}^{n'}$  and compute  $\boldsymbol{\beta} = ((\text{Neg}^i(\boldsymbol{\alpha}), \text{sk}))_{i=0}^{n'-1} + \mathbf{e} \in \mathbb{Z}_{q'}^{n'}$  with noise  $\mathbf{e} \leftarrow \chi_e^{n'}$ . Finally, the algorithm returns  $(\text{sk}, \text{pk})$ , where  $\text{pk} = (\boldsymbol{\alpha}, \boldsymbol{\beta}) \in \mathbb{Z}_{q'}^{n'} \times \mathbb{Z}_{q'}^{n'}$ .*
- $\text{sig} \leftarrow \text{PS.Signal}(\text{pk})$ : *The signaling algorithm takes a public key  $\text{pk} = (\boldsymbol{\alpha}, \boldsymbol{\beta}) \in \mathbb{Z}_{q'}^{n'} \times \mathbb{Z}_{q'}^{n'}$  and samples noises  $\mathbf{e}_0 \leftarrow \chi_s$ ,  $\mathbf{e}_1 \leftarrow \chi_e^{n'}$ , and  $\mathbf{e}_2 \leftarrow \chi_e^\ell$ . Then, it computes  $\mathbf{a} = (\langle \text{Neg}^i(\boldsymbol{\alpha}), \mathbf{e}_0 \rangle)_{i=0}^{n'-1} + \mathbf{e}_1$  and  $\mathbf{b} = (\langle \text{Neg}^i(\boldsymbol{\beta}), \mathbf{e}_0 \rangle)_{i=0}^{\ell-1} + \mathbf{e}_2$ . Finally, the algorithm returns  $\text{sig} = (\mathbf{a}, \mathbf{b}) \in \mathbb{Z}_{q'}^{n'} \times \mathbb{Z}_{q'}^\ell$ .*

- $0/1 \leftarrow \text{PS.Detect}_{\text{sk}}(\text{sig})$ : *The detection algorithm takes  $\text{sig} = (\mathbf{a}, \mathbf{b}) \in \mathbb{Z}_{q'}^{n'} \times \mathbb{Z}_{q'}^\ell$  and computes  $\mathbf{d} = \mathbf{b} - (\langle \text{Neg}^i(\mathbf{a}), \text{sk} \rangle)_{i=0}^{\ell-1} \in \mathbb{Z}_{q'}^\ell$ . Then, it computes  $c_i$ 's by checking  $\mathbf{d} = (d_1, \dots, d_\ell)$  with the range parameter  $r$ , as follows.*

$$c_i = \begin{cases} 0, & \text{if } d_i \in [-r, r] \\ 1, & \text{otherwise} \end{cases}$$

Finally, the algorithm returns  $\prod_{i=1}^{\ell} (1 - c_i) \in \{0, 1\}$ .

**Remark 2** (Security of Construction 1). *We sketch the security of Construction 1 with respect to Def. 7. For more detailed proofs, refer to [LTW24b].*

- **Completeness:** *The range parameter  $r$  can be set to guarantee completeness, following standard analysis in LWE-based encryptions to guarantee correctness.*
- **Soundness:** *The repetition parameter  $\ell$  can be set to guarantee soundness as follows: By the hardness of RLWE over  $\mathcal{R}_{q'}$ ,  $\text{sig} \leftarrow \text{PS.Signal}(\text{pk})$  is indistinguishable from a uniform random sample, in the absence of access to the secret key  $\text{sk}$ . Therefore, the following holds when  $\text{sk}'$  is sampled independently of  $\text{sk}$ .*

$$\Pr[\text{PS.Detect}_{\text{sk}'}(\text{sig}) = 1] \leq \left(\frac{2r+1}{q'}\right)^\ell + \text{negl}(\lambda)$$

- **Receiver-Privacy:** *The receiver-privacy directly follows from the hardness of RLWE over  $\mathcal{R}_{q'}$ .*

## 3.2 Setup

Given homomorphic encryption scheme BFV (Sec. 2.2.1) and private signaling scheme PS (Sec. 3.1.2), (ParamGen, KeyGen, Signal) of PerfOMR can be described as follows.

- $\text{pp} \leftarrow \text{ParamGen}(1^\lambda)$ : Run  $\text{pp}_{\text{PS}} \leftarrow \text{PS.ParamGen}(1^\lambda)$  and return  $\text{pp} = (\text{pp}_{\text{PS}}, 1^\lambda)$ .
- $(\text{sk}, \text{pk}_{\text{sig}}, \text{pk}_{\text{det}}) \leftarrow \text{KeyGen}(\text{pp})$ : Run  $(\text{sk}_{\text{PS}}, \text{pk}_{\text{PS}}) \leftarrow \text{PS.KeyGen}(\text{pp}_{\text{PS}})$  and  $(\text{sk}_{\text{BFV}}, \text{pk}_{\text{BFV}}) \leftarrow \text{BFV.KeyGen}(1^\lambda)$ . Then, compute  $\text{ct}_{\text{sk}} \leftarrow \text{BFV.Enc}(\text{sk}_{\text{PS}})$ . Finally, return  $\text{sk} = (\text{sk}_{\text{PS}}, \text{sk}_{\text{BFV}})$ ,  $\text{pk}_{\text{sig}} = \text{pk}_{\text{PS}}$ , and  $\text{pk}_{\text{det}} = (\text{pk}_{\text{BFV}}, \text{ct}_{\text{sk}})$ .
- $\text{sig} \leftarrow \text{Signal}(\text{pk}_{\text{sig}})$ : Return  $\text{sig} \leftarrow \text{PS.Signal}(\text{pk}_{\text{PS}})$ .

## 3.3 Detection

For clarity, we describe the Digest algorithm of PerfOMR in a modular manner, divided into (1) a *detection* phase and (2) a *compression* phase. In this section, we begin by illustrating Detect (Alg 6) for the detection phase, followed by an explanation of the compression phase in Sec. 3.4.

The goal of Detect is to homomorphically compute the pertinent indices from the bulletin board BB using the detection key  $\text{pk}_{\text{det}}$ . The output of Detect is ciphertexts encrypting a sparse *binary* vector, where the ones indicate the pertinent indices, referred to as the *Pertinency Vector*, PV.<sup>12</sup> At a high level, Detect can be viewed as a *batched homomorphic detection* of signals, consisting of two main steps: (i) *Affine Transform* and (ii) *Range Check*.

**Affine Transform.** The goal of Affine Transform (Step 1 of Alg. 6) is to homomorphically compute *linear* parts of multiple  $\text{PS.Detect}_{\text{skps}}(\text{sig}_i)$  in parallel. More precisely, we want to compute  $\mathbf{d} = \mathbf{b}_i - (\text{Neg}^j(\mathbf{a}_i), \text{sk})_{j=0}^{\ell-1} \in \mathbb{Z}_q^\ell$  where  $\text{sig}_i = (\mathbf{a}_i, \mathbf{b}_i) \in \mathbb{Z}_q^{n'} \times \mathbb{Z}_q^\ell$  for  $1 \leq i \leq N$ . Notice that we can stack up the inner products  $\langle \text{Neg}^j(\mathbf{a}_i), \text{sk} \rangle$  into a matrix-vector multiplication. In this way, we can leverage homomorphic MatMul algorithms (Sec. 2.3) with respect to  $\text{ct}_{\text{sk}}$ . We will explore this in greater detail in Sec. 4.4.

**Range Check.** The goal of Range Check (Step 2 of Alg. 6) is to homomorphically compute the remaining *non-linear* parts of multiple  $\text{PS.Detect}_{\text{skps}}(\text{sig}_i)$ . More precisely, we want to compute  $\prod_{i=1}^r (1 - c_i) \in \{0, 1\}$  where  $r$  is the range parameter and

$$c_i = \begin{cases} 0, & \text{if } d_i \in [-r, r] \\ 1, & \text{otherwise,} \end{cases}$$

for  $\mathbf{d} = (d_1, \dots, d_\ell)$  which is the output from Affine Transform. The computation of  $c_i$ 's can be done homomorphically by evaluating the following function on  $d_i$ 's.

$$f(X) = \begin{cases} 0, & \text{if } -r \leq X \leq r \\ 1, & \text{otherwise} \end{cases} \quad (4)$$

In PerfOMR [LTW24b],  $f(X)$  is homomorphically evaluated by first computing  $g(X) = \prod_{i=0}^r (X^2 - i^2)$  and then  $h(X) = X^{q-1}$ , rather than directly evaluating the polynomial interpolation of  $f(X)$  as done in [LT22]. Notice that  $g(X)$  maps the input to 0 if and only if it lies within the range  $[-r, r]$ , and  $h(X)$  maps non-zero inputs to 1 and zero to 0 by Fermat's little theorem. By decomposing the evaluation of  $f$  into  $g$  and  $h$ , PerfOMR significantly reduced the number of homomorphic multiplications required, while only slightly increasing the circuit depth, resulting in an overall improvement in computational cost.

## 3.4 Compression

In the compression phase, pertinent payloads are compressed into a digest using the encrypted pertinency vector PV.

<sup>12</sup>Previous works [LT22, LTW24b] have referred to this detection procedure as *ClueToPackedPV*.

---

### Algorithm 6 Detect

---

**Input:**  $\text{pk}_{\text{det}}, \text{BB} = (x_i, \text{sig}_i)_{i=1}^N \quad \triangleright N = s \cdot n$   
**Step 0:** Signal Preprocessing  $\triangleright \text{sig}_i = (\mathbf{a}_i, \mathbf{b}_i) \in \mathbb{Z}_q^{n'} \times \mathbb{Z}_q^\ell$   
**for**  $j = 0, \dots, \ell - 1$  **do**  
 $\mathbf{A}_j \leftarrow [\text{Neg}^j(\mathbf{a}_1) \dots \text{Neg}^j(\mathbf{a}_N)]^T \quad \triangleright \mathbf{A}_j \in \mathbb{Z}_q^{N \times n'}$   
 $\boldsymbol{\beta}_j \leftarrow (b_{1,j}, \dots, b_{N,j}) \quad \triangleright \mathbf{b}_i = (b_{i,j})_{j=0}^{\ell-1}$   
**Step 1:** Affine Transform  
**for**  $j = 0, \dots, \ell - 1$  **do**  
 $(\text{ct}_{i,j})_{i=0}^{s-1} \leftarrow \boldsymbol{\beta}_j - \text{MatMul}(\mathbf{A}_j, \text{ct}_{\text{sk}}) \quad \triangleright \text{see Sec. 4.4}$   
**Step 2:** Range Check  $\triangleright g(X) = X \cdot \prod_{i=1}^r (X^2 - i^2),$   
 $h(X) = X^{p-1}$   
**for**  $i = 0, \dots, s - 1$  **do**  
**for**  $j = 0, \dots, \ell - 1$  **do**  
 $\text{ct}_{i,j} \leftarrow \text{BFV.Eval}_{\text{pk}_{\text{BFV}}}(g(X), \text{ct}_{i,j})$   
 $\text{ct}_{i,j} \leftarrow \text{BFV.Eval}_{\text{pk}_{\text{BFV}}}(h(X), \text{ct}_{i,j})$   
 $\text{ct}_i \leftarrow \prod_{j=0}^{\ell-1} (1 - \text{ct}_{i,j})$   
**return**  $(\text{ct}_i)_{i=0}^{s-1}$

---

Roughly speaking, PerfOMR compresses both PV and pertinent payloads; during decoding, PV is decoded first and then used as a *hint* to decode the compressed payloads. Below is a brief overview of PerfOMR's compression phase. For a more detailed discussion, refer to [LT22, LTW24b].

#### 3.4.1 Unpacking the Pertinency Vector

In PerfOMR, before compressing the pertinency vector PV, each slot of PV are *unpacked* into separate ciphertexts. This process, referred to as PVUnpack in previous works [LT22, LTW24b], is necessary for them because they compress PV by *hashing* its components into *buckets*. (For more details, see Sec. 3.4.2.) PVUnpack incorporates so-called SlotToCoeff [GHS12a, LW23] and *oblivious expansion* [ACLS18] as subroutines, both of which are quite costly. As a result, PVUnpack becomes a bottleneck in the entire digesting process.

#### 3.4.2 Hashing-based Index Compression

PerfOMR compresses the unpacked pertinency vector  $\text{PVUnpack}(\text{PV}) = \{\text{ct}_i\}_{i=1}^N$  via *hashing*. Each index  $1 \leq i \leq N$  is randomly assigned to a *bucket*  $Y_j$  for  $1 \leq j \leq m$ . The pertinent vector is compressed by computing  $\text{Ctr}_j \leftarrow \sum_{i \in Y_j} \text{ct}_i$  and  $\text{Acc}_j \leftarrow \sum_{i \in Y_j} i \cdot \text{ct}_i$ . Notice that  $\text{Ctr}_j$  indicates the number of pertinent indices that were hashed into  $Y_j$ , revealing whether a collision occurred. If  $\text{Ctr}_j$  encrypts the value 1, the corresponding value in  $\text{Acc}_j$  is the pertinent index; otherwise, a collision has occurred and happened, and decoding will fail. To prevent such failures, PerfOMR sets the number of buck-

ets  $m$  to be sufficiently larger than  $k$ , the upper bound on the number of pertinent indices, and repeats the process multiple times.

### 3.4.3 Payload Compression with Random Linear Codes

PerfOMR then compresses pertinent payloads using random linear codes. The core idea is to compress the vector of pertinent payloads (i.e., the payloads masked by the pertinency vector) by multiplying it with a random  $K \times N$  matrix. By setting  $K$  sufficiently larger than  $k$  – the upper bound on the number of pertinent indices – it becomes possible to solve the system of equations for the pertinent payloads via Gaussian elimination with high probability, given the pertinent indices.

### 3.4.4 Bundling: Runtime/Digest-Size Trade-off

To reduce the computational overhead of PVUnpack, PerfOMR processes multiple indices simultaneously via *bundling*. Before applying PVUnpack, PerfOMR combines  $v$  indices of the pertinency vector by homomorphically adding them into a single ciphertext. While this reduces the runtime by a factor of  $v$ , it introduces a trade-off: now the receiver only knows that a bundle contains a pertinent index, requiring retrieval of all  $v$  indices and their corresponding payloads. This approach balances the reduction in PVUnpack runtime against an increase in the final digest size.

## 3.5 Optimizations

**Modulus-Switching.** *Modulus-Switching* [BGV12] is a technique of switching the ciphertext modulus of a given HE ciphertext into a smaller one. While modulus-switching is not strictly necessary in the BFV scheme [Bra12, FV12] (unlike the BGV scheme [BGV12]), it can be highly beneficial for performance optimizations. In BFV, choosing a large enough ciphertext modulus is crucial to support the desired computation, as the modulus determines the *computational budget* of a ciphertext. However, a larger modulus increases the ciphertext size and degrades the overall performance of homomorphic operations. Modulus-switching addresses this issue by reducing the modulus as the budget is consumed, preserving the remaining budget. PerfOMR leverages modulus-switching repeatedly to keep the ciphertext modulus as small as possible without compromising the budget, enhancing performance.

**Rotation Key Management.** The BFV scheme can natively support all types of rotation. However, each requires its own rotation key and storing all these keys can be quite demanding in practice. Therefore, trade-offs between storage and computational time are often advantageous. It is also important to note that the size of rotation keys scales with the size of the ciphertext modulus. However, if the rotation keys will be used only after modulus-switching, generating them in a

smaller modulus suffices. For this reason, PerfOMR reduces the number of rotation keys more aggressively for the ones used at higher ciphertext modulus levels. The approach taken by PerfOMR will be further detailed in Sec. 4.4.

## 3.6 Security

We sketch the security of PerfOMR with respect to Def. 3. For detailed proofs, refer to [LTW24b].

- **Correctness:** The correctness of PerfOMR is implied by the completeness and soundness of the underlying private signaling scheme (Rmk. 2) and the correctness of the BFV scheme (Sec. 2.2.1).
- **Receiver-Privacy:** The receiver-privacy of PerfOMR is implied by the receiver-privacy of the underlying private signaling scheme (Rmk. 2) and the semantic security of the BFV scheme (Sec. 2.2.1).

## 4 Our Scheme: SophOMR

In this section, we present our new OMR scheme, SophOMR. While we largely follow the approach of PerfOMR [LTW24b], the key distinction lies in the compression phase (Sec. 3.4), where we extend and apply the SIMD-aware homomorphic compression scheme of [CLPY24] (Sec. 2.4). A detailed explanation of our homomorphic compression mechanism is provided in Sec. 4.1. Additionally, we propose several optimizations: ring-switching for digest-size reduction (Sec. 4.2), hybrid key-switching optimization for key-size reduction (Sec. 4.3), and BSGS-style Affine Transform step (Sec. 3.3) for run-time reduction (Sec. 4.4). Additionally, we briefly discuss the OMD variant of SophOMR (Sec. 4.5).

### 4.1 SIMD-Aware Digest Compression

As discussed in Sec. 3.4, PerfOMR [LTW24b] cannot fully exploit the homomorphic SIMD structure (Sec. 2.2.1) during the compression step due to its hashing-based approach (Sec. 3.4.2). As a result, it requires the costly PVUnpack process (Sec. 3.4.1), which unpacks each slot of the packed PV ciphertext into separate ciphertexts, creating a bottleneck for the entire Digest.

In this regard, we adopt a *SIMD-aware* homomorphic compression scheme from [CLPY24] (Sec. 2.4), which removes the need for the heavy PVUnpack step, significantly improving the runtime of the compression step.<sup>13</sup> However, [CLPY24] only considered payloads of the size of a single HE slot, which is unrealistic for many scenarios. For example, in prior works on OMR [LT22, LTW24b], each payload is

<sup>13</sup>Another advantage of our compression scheme is that, unlike PerfOMR’s, it is *deterministic* and does not introduce any additional failure probability.

assumed to be 612 bytes, the typical size of Zcash transactions [HBHW], meaning that the payload occupies multiple HE slots. Thus, we need to extend the approach of [CLPY24] to accommodate these typical OMR settings.

**Naive Approach.** To be precise, let each payload occupy  $t$  number of HE slots (e.g.,  $t \approx 300$  for typical HE parameters in the example of Zcash transactions). Let  $\text{PL}_i$  denote the length- $N$  vector containing the  $i$ -th slot of payloads. Directly applying the homomorphic compression of [CLPY24] (Sec. 2.4) to the compression step of OMR (Sec. 3.4), one would first compress the encrypted pertinency vector PV obtained from the detection step (Sec. 3.3). This involves homomorphically computing  $\mathbf{C} \cdot \text{PV}$  via MatMul (Alg. 1), where  $\mathbf{C}$  is defined in Eq. 3. Then, one would compress *masked* payloads  $\text{PL}_i \odot \text{PV}$  by computing  $\mathbf{C} \cdot (\text{PL}_i \odot \text{PV})$  again with MatMul, for  $1 \leq i \leq t$ . This naive approach would require a total  $(t + 1)$  number of calls to MatMul with respect to the  $k \times N$  matrix  $\mathbf{C}$ . While this method is already fairly efficient, it can be optimized much further, as described below.

**Stacking Matrices.** First, we note a simple yet powerful fact that enables our optimization. Let  $\text{diag}(\text{PL}_i)$  denote the  $N \times N$  diagonal matrix whose diagonal is  $\text{PL}_i$  and let  $\mathbf{D}_i$  denote the  $k \times N$  matrix  $\mathbf{C} \cdot \text{diag}(\text{PL}_i)$ . Then, the following holds.

$$\mathbf{C} \cdot (\text{PL}_i \odot \text{PV}) = \mathbf{D}_i \cdot \text{PV}$$

This allows the detector to precompute  $\mathbf{D}_i = \mathbf{C} \cdot \text{diag}(\text{PL}_i)$  and then compute  $\mathbf{D}_i \cdot \text{PV}$  homomorphically. Note that this preprocessing is almost free, as both  $\mathbf{C}$  and  $\text{diag}(\text{PL}_i)$  are *plaintexts* with regard to the underlying HE scheme. Furthermore, this approach consumes less *computational budget* (Sec. 3.5) as the homomorphic computation of  $\text{PL}_i \odot \text{PV}$  is no longer performed before MatMul.

But much more importantly, this precomputation allows us to *stack* multiple MatMul into a single MatMul as follows.

$$[\mathbf{C}^T | \mathbf{D}_1^T | \dots | \mathbf{D}_t^T]^T \cdot \text{PV}$$

To see why this is useful at all, recall that the number of homomorphic rotations required in  $k \times N$  MatMul is  $O(\sqrt{k})$ . (See Tab. 1.) The stacking method significantly reduces the total number of rotations required for compression: the naive approach, which involves  $t + 1$  calls to  $k \times N$  MatMul, requires  $O((t + 1) \cdot \sqrt{k})$  rotations. In contrast, the stacking approach, which involves a single call to  $((t + 1) \cdot k) \times N$  MatMul, only requires  $O(\sqrt{(t + 1) \cdot k})$  rotations. This results in roughly a  $\sqrt{t + 1}$  reduction in the total number of rotations.

**Implementing MatMul.** We use MatMul described in Alg. 1 (Sec. 2.3) for digest compression. In particular, we employ the Horner-style [HS18] MatMul, which significantly reduces the number of required rotation keys. Although this

method is incompatible with the so-called *hoisting* technique [HS18], which optimizes runtime, we prioritize Horner-style MatMul. We believe that in scenarios where a detector handles multiple clients, minimizing the key size per client is much more important than the runtime improvements offered by hoisting.

## 4.2 Further Compression with Ring-Switching

Similar to modulus-switching (Sec. 3.5), there is a technique called *ring-switching* [BGV12, GHPS12, GHPS13], which switches the ring degree  $n$  of a given HE ciphertext into a smaller one. This can be done while preserving the same security level after some modulus-switching.<sup>14</sup> In our case, ring-switching is useful for further reducing the digest size. Specifically, the ciphertext size scales linearly with the ring dimension  $n$ , which corresponds to the number of slots in the BFV scheme (Sec. 2.2.1). This implies that if the final result occupies fewer slots than  $n$ , the unused slots are being wasted. For instance, after homomorphic compression (Sec. 4.1), only  $(t + 1) \cdot k$  slots are occupied, which is typically smaller than  $n$  under standard OMR parameters. Ring-switching addresses this inefficiency by reducing the ring dimension and making the digest more compact.<sup>15</sup> In SophOMR, ring-switching is employed as the final step of Digest to further compress the digest. We note that ring-switching is a relatively inexpensive procedure, whose cost is dominated by a single key-switching operation. For further details, please refer to [GHPS13].

## 4.3 Key Size Reduction through Hybrid Key-Switching Optimization

The *key-switching* technique is a fundamental element of homomorphic multiplications and rotations. There are two main approaches to key-switching. The first, known as the *bit-decomposition* approach [BV11, BEHZ16] optimizes the *computational budget* (Sec. 3.5). In contrast, the *special-modulus* approach [GHS12b, CHK<sup>+</sup>19] prioritize key-size optimization. In practice, a hybrid approach that combines these two methods has been found to be most favorable [GHS12b, HK20, KPZ21]. One downside of PerfOMR is its large key-size, which stems from its reliance on one extreme of this hybrid approach. In SophOMR, we strike a balance between these two extremes, significantly reducing the key-size.

## 4.4 Affine Transform with BSGS-MatMul

We optimize the runtime of Affine Transform step in the Detect procedure (Sec. 3.3) by adopting the BSGS-style Mat-

<sup>14</sup>For the same level of security, a larger ring dimension is required when using a larger ciphertext modulus.

<sup>15</sup>Ring-switching is especially effective in the context of OMD (Sec. 2.5.2), where the digest contains only the pertinent indices, excluding payloads, requiring significantly fewer slots. See Sec. 4.5.

Mul algorithm, which uses two rotation keys. Thanks to hybrid key-switching (Sec. 4.3), this optimization is achieved while maintaining a smaller key-size compared to PerfOMR’s approach of using a single rotation key. (See Sec. 5.)

Recall that the Affine Transform step involves homomorphic matrix multiplications between *tall* plaintext matrices  $\mathbf{A}_j$ ’s and the ciphertext vector  $\text{ct}_{\text{sk}}$  (Sec. 3.3, Alg. 6). For clarity, in this section, we consider the goal of homomorphically multiplying a tall  $N \times k$  plaintext matrix  $\mathbf{M}$  (representing  $\mathbf{A}_j$ ) with a ciphertext  $\text{ct}(\mathbf{v})$  (representing  $\text{ct}_{\text{sk}}$ ) that encrypts a vector  $\mathbf{v} \in \mathbb{Z}_p^k$  via sparse packing (Sec. 2.2.2).

#### 4.4.1 PerfOMR’s Approach: MatMul with Single Key

Note that the Affine Transform (Sec. 3.3) is the very first step of digesting and occurs before any modulus-switching (Sec. 3.5). As a result, the rotation keys used in this step must be generated at the largest modulus, making them relatively large. In this regard, PerfOMR [LTW24b] aggressively minimizes the number of rotation keys required during the Affine Transform, ultimately using only one. The approach skips BSGS-style optimization (Sec. 2.3.2) and applies Eq. 2 directly. Alg. 7 outlines PerfOMR’s method, with its costs summarized in Tab. 2.

---

#### Algorithm 7 MatMul w/ Single Rotation Key

---

**Input:**  $\mathbf{M} \in \mathbb{Z}_p^{N \times k}$ ,  $\text{ct}(\mathbf{v})$   $\triangleright k < n, N = s \cdot n, \mathbf{v} \in \mathbb{Z}_p^k$

**Step 0:** Parameter Setting

$$[\mathbf{M}_0^T | \dots | \mathbf{M}_{s-1}^T] \leftarrow \mathbf{M}^T$$

**Step 1:** Rotation Preprocessing

$$\begin{aligned} \text{ct}_0 &\leftarrow \text{ct}(\mathbf{v}) \\ \text{for } i = 1, \dots, k-1 \text{ do} \\ \text{ct}_i &\leftarrow \text{Rot}_{\text{row}}^1(\text{ct}_{i-1}) \end{aligned}$$

**Step 2:** Ptxt-Mult & Sum

$$\text{for } j = 0, \dots, s-1 \text{ do} \\ \text{ct}_{\text{out},j} \leftarrow \sum_{i=0}^{k-1} \text{diag}_i(\mathbf{M}_j) \odot \text{ct}_i$$

**return**  $(\text{ct}_{\text{out},j})_{j=0}^{s-1}$

---

#### 4.4.2 Our Approach: MatMul with Two Keys

In SophOMR, thanks to hybrid key-switching (Sec. 4.3), the size of the rotation key generated at the largest modulus is substantially smaller than in PerfOMR. As a result, even with the BSGS-style optimization and the use of two keys, SophOMR’s total key-size remains smaller than PerfOMR’s. The adoption of BSGS reduces the number of required homomorphic rotations from  $O(k)$  to  $O(\sqrt{k})$ , enhancing runtime performance.

Alg. 8 describes the BSGS-MatMul algorithm used for Affine transform in SophOMR. Although it is simply the tall

matrix version of Alg. 1, we include it here for completeness. The costs for Alg. 8 are summarized in Tab. 2. Note that the number of rotations is minimized when  $\tilde{b} \approx \sqrt{s \cdot k}$ .

---

#### Algorithm 8 MatMul (Tall)

---

**Input:**  $\mathbf{M} \in \mathbb{Z}_p^{N \times k}$ ,  $\text{ct}(\mathbf{v})$   $\triangleright k < n, N = s \cdot n, \mathbf{v} \in \mathbb{Z}_p^k$

**Step 0:** Parameter Setting

$$[\mathbf{M}_0^T | \dots | \mathbf{M}_{s-1}^T] \leftarrow \mathbf{M}^T \text{ and } k = \tilde{g} \cdot \tilde{b}$$

**Step 1:** Plaintext Packing

$$\text{for } (i, j) \in [s] \times [k] \text{ do} \\ \mathbf{m}_{i,j} \leftarrow \text{Rot}_{\text{row}}^{\tilde{g}\tilde{b}}(\text{diag}_j(\mathbf{M}_i)) \text{ for } g = \lfloor j/\tilde{b} \rfloor$$

**Step 2:** Baby-Step-Giant-Step (BSGS)

$$\begin{aligned} \text{ct}_0 &\leftarrow \text{ct}(\mathbf{v}) && \triangleright \text{Baby-Step} \\ \text{for } b = 1, \dots, \tilde{b}-1 \text{ do} \\ \text{ct}_b &\leftarrow \text{Rot}_{\text{row}}^1(\text{ct}_{b-1}) \\ \text{for } i = 0, \dots, s-1 \text{ do} && \triangleright \text{Giant-Step} \\ \text{for } g = \tilde{g}-1, \dots, 0 \text{ do} \\ \text{ct}_{\text{sum},i} &\leftarrow \sum_{b=0}^{\tilde{b}-1} (\mathbf{m}_{i,g\tilde{b}+b} \odot \text{ct}_b) \\ \text{if } g = \tilde{g}-1 \text{ then} \\ \text{ct}_{\text{out},i} &\leftarrow \text{ct}_{\text{sum},i} \\ \text{else} \\ \text{ct}_{\text{out},i} &\leftarrow \text{Rot}_{\text{row}}^{\tilde{b}}(\text{ct}_{\text{out},i}) + \text{ct}_{\text{sum},i} \end{aligned}$$

**return**  $(\text{ct}_{\text{out},i})_{i=0}^{s-1}$

---

	#(Ptxt-Mult)	#(Rotation)	#(Rot. Key)
Alg. 7	$s \cdot k$	$k$	1
Alg. 8	$s \cdot k$	$\tilde{b} + s \cdot \tilde{g}$	2

Table 2: Cost Analysis for Alg. 7 and Alg. 8

## 4.5 OMD Variant

Our SophOMR scheme can be readily adapted into an Oblivious Message Detection (OMD) scheme (Sec. 2.5.2). Since OMD does not require payload compression, it is sufficient to compress only the encrypted pertinency vector PV. This is done by homomorphically computing  $\mathbf{C} \cdot \text{PV}$  via MatMul (Alg. 1), where  $\mathbf{C}$  is defined in Eq. 3, without stacking the matrices as described in Sec. 4.1. This leads to a significant speedup in the compression phase (Sec. 4.1). See Sec. 5 for concrete performance results. An OMR-like functionality can be achieved by combining an OMD scheme with a private information retrieval (PIR) scheme [CGKS95, KO97], as discussed in Sec. 2.5.2. Although this approach introduces an additional round of communication, it provides a generic method to implement OMR-like functionality, allowing the unique features and advantages of various PIR schemes to be utilized depending on the specific PIR scheme chosen.

## 5 Evaluation

**Methodology.** To evaluate the concrete performance of our SophOMR scheme, we implement it in a C++ library<sup>16</sup> using the OpenFHE library [BAB<sup>+</sup>22], with the NTL library [Sho] to implement Decode (Sec. 4.1 and 2.4). We then compare its performance against the implementation [LTW] of PerfOMR [LTW24b], which uses PALISADE library [PAL21] for the private signaling scheme (Sec. 3.1), the SEAL library [SEA23] for the BFV scheme (Sec. 2.2.1). Both implementations leverage the HEXL library [BKS<sup>+</sup>21] to accelerate HE libraries. All experiments are conducted on a Google Compute Cloud `n2-standard-8` instance (Intel Ice Lake) in a single-threaded mode.

**Parameters.** We choose parameters that uphold the security claims of PerfOMR’s parameter sets and comply with the Homomorphic Encryption Security Standardization effort [ACC<sup>+</sup>18, BCC<sup>+</sup>24]. All bit security is measured using the latest version of the Lattice Estimator available at the time of writing [Aib, APS15] (Commit ID: 5c25455). Detailed parameters are provided in Tab. 4, Tab. 5, and Tab. 6 of Appx. A. Note that SophOMR uses a larger BFV ring dimension of  $n = 2^{16}$ , compared to  $n = 2^{15}$  of PerfOMR, to comply with the HE Security Standard for 128-bit security (Tab. 4). We also employ a larger parameter for the underlying private signaling scheme (Sec. 3.1) that offers 128-bit security and 30-bit completeness to uphold the security claim of PerfOMR’s parameters. The parameters provided in PerfOMR only achieve 94-bit security and 18-bit completeness, leading to easily observable errors (Tab. 6).

**Experiments.** Experiments are conducted with a total number of payloads  $N = 2^{16}$  and  $2^{19}$ . The number of pertinent payloads is set to  $k = 50$ , following PerfOMR. We present the results in Tab. 3 for  $N = 2^{16}$ . Additionally, Fig. 2 illustrate the runtime breakdown of each component contributing to the Digest runtime: Preprocessing (Rmk. 4), Affine Transform (Sec. 4.4), Range Check (Sec. 3.3), and Compress (Sec. 4.1). Refer to Appx. B, for experimental results with  $N = 2^{19}$ .

**Remark 3 (On Choice of  $N$ ).** While PerfOMR conducted experiments with  $N = 2^{19}, 2^{21}, 2^{23}$ , we focus on  $N = 2^{16}, 2^{19}$ , considering  $N = 2^{16}$  as the primary setting. Our reasoning is that dividing  $N$  into smaller sets can be advantageous for latency, provided it maintains similar throughput and keeps communication costs reasonable. For example, PerfOMR uses  $N = 2^{19}$  to model Bitcoin-scale applications,<sup>17</sup> but a client may find it more useful to receive digests for  $2^{16}$  payloads every 3 hours than a single digest for  $2^{19}$  every 24 hours. This issue is even more relevant in private messaging scenarios.

<sup>16</sup><https://github.com/keewoolee/SophOMR>

<sup>17</sup>There are roughly  $2^{19}$  transactions per day on Bitcoin: [https://ycharts.com/indicators/bitcoin\\_transactions\\_per\\_day](https://ycharts.com/indicators/bitcoin_transactions_per_day)

**Remark 4 (Preprocessing).** In Fig. 2 and Fig. 3, the blue mesh represents the time that PerfOMR manually excluded from their timing results to account for the preprocessing potential of a significant portion of SlotToCoeff (Sec. 3.4.1). However, this preprocessing requires a substantial amount of RAM (at least 25GB for their parameters). If such an amount of RAM is available, SophOMR can also be accelerated by preprocessing a significant portion of MatMul in Compress (Sec. 4.1).<sup>18</sup> For a fair comparison, preprocessing time is excluded from Tab. 3 and Tab. 7, but it is provided in Fig. 2 and Fig. 3.

**Experimental Results for  $N = 2^{16}$ .** As shown in Tab. 3, SophOMR provides multiple improvements over PerfOMR. Specifically, SophOMR achieves a  $3.3\times$  speedup in Digest runtime while maintaining a reasonably small Decode runtime. The detection key size is reduced by  $1.3\times$  through hybrid key-switching optimization (Sec. 4.3). Furthermore, the digest size is reduced by  $2.2\times$ , thanks to ring-switching (Sec. 4.2) and our compact homomorphic compression scheme (Sec. 4.1); see also Rmk. 6. Note that all these improvements are achieved despite SophOMR using larger parameters than PerfOMR to meet security levels that PerfOMR does not satisfy (Appx. A).

As illustrated in Fig. 2, the improvement in Digest runtime primarily stems from our SIMD-aware homomorphic compression scheme (Sec. 4.1). Specifically, our compression achieves  $7.4\times$  speedup compared to PerfOMR, even when excluding the runtime for preprocessing (Rmk. 4). The Affine Transform step also achieves considerable speedup thanks to BSGS-style MatMul (Sec. 4.4). On the other hand, the Range Check step (Sec. 3.3) shows a slight increase in runtime, attributed to our use of larger parameters to meet security levels that PerfOMR does not satisfy (Appx. A).

**Remark 5 (On Choice of  $v$ ).** The bundling parameter  $v$  in PerfOMR represents how many payloads are concatenated and processed simultaneously (See Sec. 3.4.4). We compare only with PerfOMR using  $v = 2$  for  $N = 2^{16}$  because (i) setting  $v = 2$  incurs nearly no additional runtime or digest size costs compared to  $v = 1$  due to PerfOMR’s specific utilization BFV slots, and (ii) setting  $v > 2$  does not improve runtime further due to the fixed number of BFV slots ( $n = 2^{15}$ ).

**Remark 6 (Digest Size).** PerfOMR optimizes the digest size by setting the BFV ciphertext modulus at the lowest level to be as small as possible in the SEAL library (Sec. 3.5). However, this optimization is not applied in our SophOMR implementation due to limitations in the OpenFHE library: (i) it does not officially allow the base-level ciphertext modulus to be set independently of the RNS limbs, and (ii) its ciphertext serialization does not further reduce when the ciphertext

<sup>18</sup>If an even larger RAM is available, we can also preprocess Affine Transform (Sec. 4.4).

	Digest (s)	Decode (ms)	Detection Key (MB) <sup>19</sup>	Digest (KB) <sup>20</sup>
PerfOMR ( $v = 2$ )	534	12	297	567
SophOMR	162	14	228	263
PerfOMD ( $v = 2$ )	492	4	297	284
SophOMD	113	5	284	132

Table 3: Performance of OMR/OMD Schemes ( $N = 2^{16}$ ,  $k = 50$ )

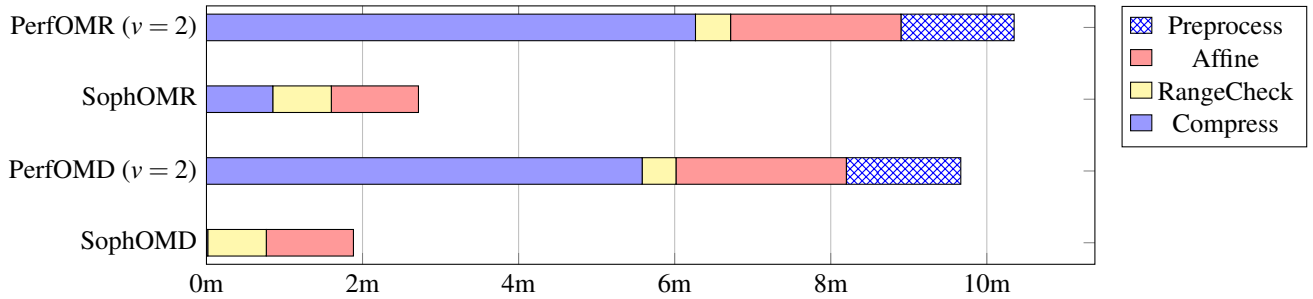


Figure 2: Runtime Breakdown for OMR/OMD Schemes ( $N = 2^{16}$ ,  $k = 50$ )

modulus is below 64-bit. If we apply the same optimization as in PerfOMR, SophOMR’s digest size can be reduced by approximately an additional factor of  $2\times$ .

parameter  $v$ .

**OMD Variants.** We also present the performance of SophOMD, the OMD variant of SophOMR (see Sec. 4.5). As there is no existing OMD implementation corresponding to PerfOMR, we modified the PerfOMR implementation to run only the relevant components and refer to this version as PerfOMD. The OMD experiments use the same parameter sets as in the OMR experiments.

As shown in Tab. 3, SophOMD also provides multiple improvements over PerfOMD, similar to the OMR setting. Specifically, SophOMD achieves a  $4.4\times$  speedup in Digest runtime. As illustrated in Fig. 2, this improvement primarily stems from our SIMD-aware homomorphic compression scheme (Sec. 4.1), allowing compression in SophOMD to complete in only one second – a  $335\times$  speedup compared to PerfOMD. Notice that the performance gain is even greater than in the OMR setting, as we only compress the indices and not the payloads (Sec. 4.5). On the other hand, the detector key size has slightly increased (though still smaller than PerfOMR) due to the use of additional rotation keys from the reduced number of rows in MatMul (Tab. 1). The bundling technique (Sec. 3.4.4) can still be applied to PerfOMD; however, achieving OMR-like functionality (Sec. 4.5) requires performing PIR (Private Information Retrieval) on the *bundles*, causing PIR performance to depend on the bundling

<sup>19</sup>The detector key size is larger than the values reported in [LTW24b] because they did not account for the size of relinearization keys.

<sup>20</sup>Our digest size can be further reduced, as explained in Rmk. 6.

## References

- [ACC<sup>+</sup>18] Martin Albrecht, Melissa Chase, Hao Chen, Jintai Ding, Shafi Goldwasser, Sergey Gorbunov, Shai Halevi, Jeffrey Hoffstein, Kim Laine, Kristin Lauter, Satya Lokam, Daniele Micciancio, Dustin Moody, Travis Morrison, Amit Sahai, and Vinod Vaikuntanathan. Homomorphic encryption security standard. Technical report, HomomorphicEncryption.org, Toronto, Canada, November 2018.
- [ACLS18] Sebastian Angel, Hao Chen, Kim Laine, and Srinath 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.
- [Alb] Martin Albrecht. Security Estimates for Lattice Problems. <https://github.com/malb/lattice-estimator>.
- [APS15] Martin R. Albrecht, Rachel Player, and Sam Scott. On the concrete hardness of learning with errors. *J. Math. Cryptol.*, 9(3):169–203, 2015.
- [AS16] Sebastian Angel and Srinath T. V. Setty. Unobservable communication over fully untrusted infrastructure. In Kimberly Keeton and Timothy Roscoe, editors, *12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016, Savannah, GA, USA, November 2-4, 2016*, pages 551–569. USENIX Association, 2016.
- [Azt] Aztec. <https://aztec.network/>.
- [BAB<sup>+</sup>22] Ahmad Al Badawi, Andreea Alexandru, Jack Bates, Flavio Bergamaschi, David Bruce Cousins, Saroja Erabelli, Nicholas Genise, Shai Halevi, Hamish Hunt, Andrey Kim, Yongwoo Lee, Zeyu Liu, Daniele Micciancio, Carlo Pascoe, Yuriy Polyakov, Ian Quah, Saraswathy R.V., Kurt Rohloff, Jonathan Saylor, Dmitriy Sponitsky, Matthew Triplett, Vinod Vaikuntanathan, and Vincent Zucca. OpenFHE: Open-source fully homomorphic encryption library. Cryptology ePrint Archive, Paper 2022/915, 2022. <https://eprint.iacr.org/2022/915>.
- [BCC<sup>+</sup>24] Jean-Philippe Bossuat, Rosario Cammarota, Iliaria Chillotti, Benjamin R. Curtis, Wei Dai, Huijing Gong, Erin Hales, Duhyeong Kim, Bryan Kumara, Changmin Lee, Xianhui Lu, Carsten Maple, Alberto Pedrouzo-Ulloa, Rachel Player, Yuriy Polyakov, Luis Antonio Ruiz Lopez, Yongsoo Song, and Donggeon Yhee. Security guidelines for implementing homomorphic encryption. Cryptology ePrint Archive, Paper 2024/463, 2024.
- [BCG<sup>+</sup>14] Eli Ben-Sasson, Alessandro Chiesa, Christina Garman, Matthew Green, Ian Miers, Eran Tromer, and Madars Virza. Zerocash: Decentralized anonymous payments from bitcoin. In *2014 IEEE Symposium on Security and Privacy*, pages 459–474. IEEE Computer Society Press, May 2014.
- [BEHZ16] Jean-Claude Bajard, Julien Eynard, M. Anwar Hasan, and Vincent Zucca. A full RNS variant of FV like somewhat homomorphic encryption schemes. In Roberto Avanzi and Howard M. Heys, editors, *SAC 2016*, volume 10532 of *LNCS*, pages 423–442. Springer, Cham, August 2016.
- [BGV12] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (Leveled) fully homomorphic encryption without bootstrapping. In Shafi Goldwasser, editor, *ITCS 2012*, pages 309–325. ACM, January 2012.
- [BKS<sup>+</sup>21] Fabian Boemer, Sejun Kim, Gelila Seifu, Filipe DM de Souza, Vinodh Gopal, et al. Intel HEXL (release 1.2). <https://github.com/intel/hexl>, 2021.
- [BLMG21] Gabrielle Beck, Julia Len, Ian Miers, and Matthew Green. Fuzzy message detection. In Giovanni Vigna and Elaine Shi, editors, *ACM CCS 2021*, pages 1507–1528. ACM Press, November 2021.
- [BMS01] Adam Back, Ulf Möller, and Anton Stiglic. Traffic analysis attacks and trade-offs in anonymity providing systems. In Ira S. Moskowitz, editor, *Information Hiding, 4th International Workshop, IHW 2001, Pittsburgh, PA, USA, April 25-27, 2001, Proceedings*, volume 2137 of *Lecture Notes in Computer Science*, pages 245–257. Springer, 2001.
- [BPSY24] Alexander Bienstock, Sarvar Patel, Joon Young Seo, and Kevin Yeo. Batch PIR and labeled PSI with oblivious ciphertext compression. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 5949–5966, Philadelphia, PA, August 2024. USENIX Association.



- [Bra12] Zvika Brakerski. Fully homomorphic encryption without modulus switching from classical GapSVP. In Reihaneh Safavi-Naini and Ran Canetti, editors, *CRYPTO 2012*, volume 7417 of *LNCS*, pages 868–886. Springer, Berlin, Heidelberg, August 2012.
- [But20a] Vitalik Buterin. Exploring fully homomorphic encryption. <https://vitalik.eth.limo/general/2020/07/20/homomorphic.html>, 2020.
- [But20b] Vitalik Buterin. Open problem: improving stealth addresses. <https://ethresear.ch/t/open-problem-improving-stealth-addresses/7438>, 2020.
- [BV11] Zvika Brakerski and Vinod Vaikuntanathan. Fully homomorphic encryption from ring-LWE and security for key dependent messages. In Phillip Rogaway, editor, *CRYPTO 2011*, volume 6841 of *LNCS*, pages 505–524. Springer, Berlin, Heidelberg, August 2011.
- [CBM15] Henry Corrigan-Gibbs, Dan Boneh, and David Mazières. Riposte: An anonymous messaging system handling millions of users. In *2015 IEEE Symposium on Security and Privacy*, pages 321–338. IEEE Computer Society Press, May 2015.
- [CDG<sup>+</sup>21] Seung Geol Choi, Dana Dachman-Soled, S. Dov Gordon, Linsheng Liu, and Arkady Yerukhimovich. Compressed oblivious encoding for homomorphically encrypted search. In Giovanni Vigna and Elaine Shi, editors, *ACM CCS 2021*, pages 2277–2291. ACM Press, November 2021.
- [CDKS19] Hao Chen, Wei Dai, Miran Kim, and Yongsoo Song. Efficient multi-key homomorphic encryption with packed ciphertexts with application to oblivious neural network inference. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *ACM CCS 2019*, pages 395–412. ACM Press, November 2019.
- [CGGI20] Iliaria Chillotti, Nicolas Gama, Mariya Georgieva, and Malika Izabachène. TFHE: Fast fully homomorphic encryption over the torus. *Journal of Cryptology*, 33(1):34–91, January 2020.
- [CGKS95] Benny Chor, Oded Goldreich, Eyal Kushilevitz, and Madhu Sudan. Private information retrieval. In *36th FOCS*, pages 41–50. IEEE Computer Society Press, October 1995.
- [CHK<sup>+</sup>19] Jung Hee Cheon, Kyoohyung Han, Andrey Kim, Miran Kim, and Yongsoo Song. A full RNS variant of approximate homomorphic encryption. In Carlos Cid and Michael J. Jacobson Jr., editors, *SAC 2018*, volume 11349 of *LNCS*, pages 347–368. Springer, Cham, August 2019.
- [CKKS17] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yong Soo Song. Homomorphic encryption for arithmetic of approximate numbers. In Tsuyoshi Takagi and Thomas Peyrin, editors, *ASIACRYPT 2017, Part I*, volume 10624 of *LNCS*, pages 409–437. Springer, Cham, December 2017.
- [CLPY24] Jung Hee Cheon, Keewoo Lee, Jai Hyun Park, and Yongdong Yeo. SIMD-aware homomorphic compression and application to private database query, 2024. <https://arxiv.org/abs/2408.17063>.
- [FLS23] Nils Fleischhacker, Kasper Green Larsen, and Mark Simkin. How to compress encrypted data. In Carmit Hazay and Martijn Stam, editors, *EUROCRYPT 2023, Part I*, volume 14004 of *LNCS*, pages 551–577. Springer, Cham, April 2023.
- [FV12] Junfeng Fan and Frederik Vercauteren. Somewhat practical fully homomorphic encryption. *Cryptology ePrint Archive*, Report 2012/144, 2012.
- [Gen09] Craig Gentry. Fully homomorphic encryption using ideal lattices. In Michael Mitzenmacher, editor, *41st ACM STOC*, pages 169–178. ACM Press, May / June 2009.
- [GHPS12] Craig Gentry, Shai Halevi, Chris Peikert, and Nigel P. Smart. Ring switching in BGV-style homomorphic encryption. In Ivan Visconti and Roberto De Prisco, editors, *SCN 12*, volume 7485 of *LNCS*, pages 19–37. Springer, Berlin, Heidelberg, September 2012.
- [GHPS13] Craig Gentry, Shai Halevi, Chris Peikert, and Nigel P. Smart. Field switching in bgv-style homomorphic encryption. *J. Comput. Secur.*, 21(5):663–684, 2013.
- [GHS12a] Craig Gentry, Shai Halevi, and Nigel P. Smart. Better bootstrapping in fully homomorphic encryption. In Marc Fischlin, Johannes Buchmann, and Mark Manulis, editors, *PKC 2012*,

- volume 7293 of *LNCS*, pages 1–16. Springer, Berlin, Heidelberg, May 2012.
- [GHS12b] Craig Gentry, Shai Halevi, and Nigel P. Smart. Homomorphic evaluation of the AES circuit. In Reihaneh Safavi-Naini and Ran Canetti, editors, *CRYPTO 2012*, volume 7417 of *LNCS*, pages 850–867. Springer, Berlin, Heidelberg, August 2012.
- [HBHW] Daira-Emma Hopwood, Sean Bowe, Taylor Hornby, and Nathan Wilcox. Zcash Protocol Specification. <https://zips.z.cash/protocol/protocol.pdf>.
- [HK20] Kyohyung Han and Dohyeong Ki. Better bootstrapping for approximate homomorphic encryption. In Stanislaw Jarecki, editor, *CT-RSA 2020*, volume 12006 of *LNCS*, pages 364–390. Springer, Cham, February 2020.
- [HS14] Shai Halevi and Victor Shoup. Algorithms in HElib. In Juan A. Garay and Rosario Gennaro, editors, *CRYPTO 2014, Part I*, volume 8616 of *LNCS*, pages 554–571. Springer, Berlin, Heidelberg, August 2014.
- [HS18] Shai Halevi and Victor Shoup. Faster homomorphic linear transformations in HElib. In Hovav Shacham and Alexandra Boldyreva, editors, *CRYPTO 2018, Part I*, volume 10991 of *LNCS*, pages 93–120. Springer, Cham, August 2018.
- [JLM23] Sashidhar Jakkamsetti, Zeyu Liu, and Varun Madathil. Scalable private signaling. Cryptology ePrint Archive, Paper 2023/572, 2023. <https://eprint.iacr.org/2023/572>.
- [JMK24] Yanxue Jia, Varun Madathil, and Aniket Kate. HomeRun: High-efficiency oblivious message retrieval, unrestricted. Cryptology ePrint Archive, Paper 2024/188, 2024. <https://eprint.iacr.org/2024/188>.
- [JVC18] Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha Chandrakasan. GAZELLE: A low latency framework for secure neural network inference. In William Enck and Adrienne Porter Felt, editors, *USENIX Security 2018*, pages 1651–1669. USENIX Association, August 2018.
- [KO97] Eyal Kushilevitz and Rafail Ostrovsky. Replication is NOT needed: SINGLE database, computationally-private information retrieval. In *38th FOCS*, pages 364–373. IEEE Computer Society Press, October 1997.
- [KPZ21] Andrey Kim, Yuriy Polyakov, and Vincent Zucca. Revisiting homomorphic encryption schemes for finite fields. In Mehdi Tibouchi and Huaxiong Wang, editors, *ASIACRYPT 2021, Part III*, volume 13092 of *LNCS*, pages 608–639. Springer, Cham, December 2021.
- [LHH<sup>+</sup>21] Wen-jie Lu, Zhicong Huang, Cheng Hong, Yiping Ma, and Hunter Qu. PEGASUS: bridging polynomial and non-polynomial evaluations in homomorphic encryption. In *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*, pages 1057–1073. IEEE, 2021.
- [LPR10] Vadim Lyubashevsky, Chris Peikert, and Oded Regev. On ideal lattices and learning with errors over rings. In Henri Gilbert, editor, *EUROCRYPT 2010*, volume 6110 of *LNCS*, pages 1–23. Springer, Berlin, Heidelberg, May / June 2010.
- [LPR13] Vadim Lyubashevsky, Chris Peikert, and Oded Regev. A toolkit for ring-LWE cryptography. In Thomas Johansson and Phong Q. Nguyen, editors, *EUROCRYPT 2013*, volume 7881 of *LNCS*, pages 35–54. Springer, Berlin, Heidelberg, May 2013.
- [LSTW24] Zeyu Liu, Katerina Sotiraki, Eran Tromer, and Yunhao Wang. Snake-eye resistance from LWE for oblivious message retrieval and robust encryption. Cryptology ePrint Archive, Paper 2024/510, 2024. <https://eprint.iacr.org/2024/510>.
- [LT22] Zeyu Liu and Eran Tromer. Oblivious message retrieval. In Yevgeniy Dodis and Thomas Shrimpton, editors, *CRYPTO 2022, Part I*, volume 13507 of *LNCS*, pages 753–783. Springer, Cham, August 2022.
- [LTW] Zeyu Liu, Eran Tromer, and Yunhao Wang. PerfOMR: proof of concept C++ implementation for OMR (Oblivious Message Retrieval) with Reduced Communication and Computation. <https://github.com/ObliviousMessageRetrieval/ObliviousMessageRetrieval/tree/perfomr>.
- [LTW24a] Z. Liu, E. Tromer, and Y. Wang. Group oblivious message retrieval. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 118–118, Los Alamitos, CA, USA, may 2024. IEEE Computer Society.

- [LTW24b] Zeyu Liu, Eran Tromer, and Yunhao Wang. PerfOMR: Oblivious message retrieval with reduced communication and computation. In Davide Balzarotti and Wenyuan Xu, editors, *33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024*. USENIX Association, 2024.
- [LW23] Zeyu Liu and Yunhao Wang. Amortized functional bootstrapping in less than 7 ms, with  $\tilde{O}(1)$  polynomial multiplications. In Jian Guo and Ron Steinfeld, editors, *ASIACRYPT 2023, Part VI*, volume 14443 of *LNCS*, pages 101–132. Springer, Singapore, December 2023.
- [MD04] Nick Mathewson and Roger Dingledine. Practical traffic analysis: Extending and resisting statistical disclosure. In David M. Martin Jr. and Andrei Serjantov, editors, *Privacy Enhancing Technologies, 4th International Workshop, PET 2004, Toronto, Canada, May 26-28, 2004, Revised Selected Papers*, volume 3424 of *Lecture Notes in Computer Science*, pages 17–34. Springer, 2004.
- [MSS<sup>+</sup>22] Varun Madathil, Alessandra Scafuro, István András Seres, Omer Shlomovits, and Denis Varlakov. Private signaling. In Kevin R. B. Butler and Kurt Thomas, editors, *USENIX Security 2022*, pages 3309–3326. USENIX Association, August 2022.
- [Noe15] Shen Noether. Ring signature confidential transactions for monero. *Cryptology ePrint Archive*, Report 2015/1098, 2015.
- [PAL21] PALISADE Lattice Cryptography Library (release 1.11.3). <https://palisade-crypto.org/>, May 2021.
- [PVW08] Chris Peikert, Vinod Vaikuntanathan, and Brent Waters. A framework for efficient and composable oblivious transfer. In David Wagner, editor, *CRYPTO 2008*, volume 5157 of *LNCS*, pages 554–571. Springer, Berlin, Heidelberg, August 2008.
- [RAD78] Ronald L. Rivest, Len Adleman, and Michael L. Dertouzos. On data banks and privacy homomorphisms. *Foundations of Secure Computation, Academia Press*, pages 169–179, 1978.
- [Reg05] Oded Regev. On lattices, learning with errors, random linear codes, and cryptography. In Harold N. Gabow and Ronald Fagin, editors, *37th ACM STOC*, pages 84–93. ACM Press, May 2005.
- [SEA23] Microsoft SEAL (release 4.1). <https://github.com/Microsoft/SEAL>, January 2023. Microsoft Research, Redmond, WA.
- [Sho] Victor Shoup. NTL: A Library for doing Number Theory. <https://libntl.org/>.
- [SPB22] István András Seres, Balázs Pejő, and Péter Burcsi. The effect of false positives: Why fuzzy message detection leads to fuzzy privacy guarantees? In Ittay Eyal and Juan A. Garay, editors, *FC 2022*, volume 13411 of *LNCS*, pages 123–148. Springer, Cham, May 2022.
- [vdHLZZ15] Jelle van den Hooff, David Lazar, Matei Zaharia, and Nickolai Zeldovich. Vuvuzela: scalable private messaging resistant to traffic analysis. In Ethan L. Miller and Steven Hand, editors, *Proceedings of the 25th Symposium on Operating Systems Principles, SOSP 2015, Monterey, CA, USA, October 4-7, 2015*, pages 137–152. ACM, 2015.
- [WCFJ12] David Isaac Wolinsky, Henry Corrigan-Gibbs, Bryan Ford, and Aaron Johnson. Dissent in numbers: Making strong anonymity scale. In Chandu Thekkath and Amin Vahdat, editors, *10th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2012, Hollywood, CA, USA, October 8-10, 2012*, pages 179–182. USENIX Association, 2012.
- [WSDB22] Toni Wahrstätter, Matt Solomon, Ben DiFrancesco, and Vitalik Buterin. ERC-5564: Stealth Addresses. *Ethereum Improvement Proposals*, no. 5564, August 2022. Available: <https://eips.ethereum.org/EIPS/eip-5564>.

## A Parameters

The BFV parameters used in our experiments are presented in Tab. 4. For SophOMR and SophOMD, we also present the BFV parameters for ring-switching (Sec. 4.2) in Tab. 5. The private signaling parameters are presented in Tab. 6.

**BFV Parameters.** In Tab. 4 and Tab. 5,  $\log PQ$  denotes the largest modulus used in the scheme, including *special modulus* used in the hybrid key-switching (Sec. 4.3), which impacts the security level. In Tab. 4, notice that SophOMR uses a larger BFV ring dimension of  $n = 2^{16}$ , compared to  $n = 2^{15}$  of PerfOMR, to comply with the HE Security Standard for 128-bit security. In this setup, a larger plaintext modulus  $p$  is required to satisfy  $p = 1 \pmod{2n}$  to fully leverage the homomorphic SIMD structure. In Tab. 5, observe that we

further reduce the ring dimension  $n$  in SophOMD, as the final digest occupies much fewer slots than SophOMR.

	$n$	$p$	$\log Q$	$\log PQ$	$\lambda$
PerfOMR	$2^{15}$	65537	857	917	120.9
SophOMR	$2^{16}$	786433	1140	1740	128.6

Table 4: Parameters for BFV

	$n$	$\log Q$	$\log PQ$	$\lambda$
SophOMR	$2^{14}$	120	240	244.6
SophOMD	$2^{13}$	120	180	154.4

Table 5: BFV Parameters for Ring-Switching

**Private Signaling Parameters.** In Tab. 6,  $\rho$  and  $\nu$  denote the false-positive and false-negative security parameters, which are associated with soundness and completeness, respectively (Def. 7). Note that the parameters used in PerfOMR provide only 94-bit security and 18-bit completeness, falling short of their claimed 128-bit security and 30-bit completeness. When running PerfOMR for  $N = 2^{19}$  with these parameter sets, errors become readily apparent due to this lack of completeness. In SophOMR, we choose parameters that uphold the security claims of PerfOMR’s parameter sets, which negatively affect our performance. Note that our key and signal sizes are slightly larger than PerfOMR due to our larger BFV parameter, which is to comply with the HE Security Standard for 128-bit security.

## B Experimental Results for $N = 2^{19}$

In this section, we present our experimental results for  $N = 2^{19}$ . We compare the performance of SophOMR/SophOMD with PerfOMR/PerfOMD for  $\nu = 2$  and  $\nu = 16$ , representing the two extremes in the trade-off between digest size and runtime. (See Sec. 3.4.4 and Rmk. 5.) As shown in Tab. 7 and Fig. 3, SophOMR/SophOMD provides multiple improvements over PerfOMR/PerfOMD for  $N = 2^{19}$ .

**OMR Schemes.** Compared to PerfOMR with  $\nu = 2$ , SophOMR achieves a  $3.0\times$  speedup in Digest runtime and a  $2.2\times$  reduction in digest size. In contrast, compared to PerfOMR with  $\nu = 16$ , SophOMR achieves  $1.1\times$  speedup in Digest runtime and a  $10.8\times$  reduction in digest size. Additionally, note that the Decode runtime for PerfOMR increases when  $\nu = 16$ .

<sup>21</sup>For the key size, we take into account the compression of randomness into the hash seed.

**OMD Schemes.** For the OMD variants, similar trends are observed as in the  $N = 2^{16}$  setting. Specifically, our SIMD-aware homomorphic compression scheme (Sec. 4.1) enables SophOMD to complete compression for  $N = 2^{19}$  in just 2 seconds—yielding a  $1343\times$  speedup over PerfOMD with  $\nu = 2$  and a  $190\times$  speedup with  $\nu = 16$ .

	$n'$	$q'$	$\sigma$	$h$	$\ell$	$r$	$\lambda$	$\rho$	$v$	PS Key (KB) <sup>21</sup>	Signal (KB)
PerfOMR	1024	65537	0.5	32	2	19	94.0	21.4	17.6	2.2	2.2
SophOMR	1024	786433	0.5	80	2	40	128.4	26.5	30.7	2.6	2.6

Table 6: Parameters for PS

	Digest (s)	Decode (ms)	Detection Key (MB)	Digest (KB)
PerfOMR ( $v = 2$ )	3627	12	297	568
PerfOMR ( $v = 16$ )	1386	70	297	2840
SophOMR	1107	14	228	263
PerfOMD ( $v = 2$ )	3331	4	297	284
PerfOMD ( $v = 16$ )	1004	5	297	568
SophOMD	723	5	284	132

Table 7: Performance of OMR/OMD Schemes ( $N = 2^{19}$ ,  $k = 50$ )

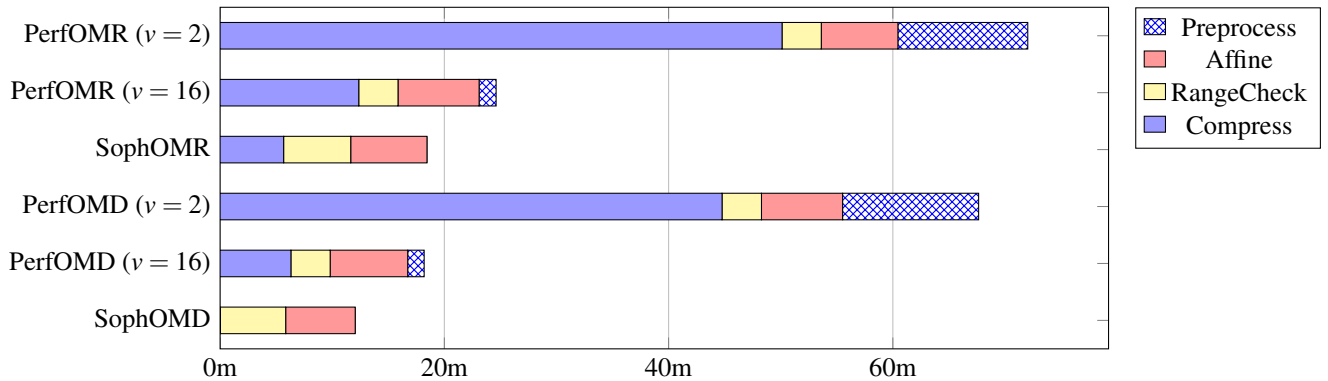


Figure 3: Runtime Breakdown for OMR/OMD Schemes ( $N = 2^{19}$ ,  $k = 50$ , single-thread)