

PerfOMR: Oblivious Message Retrieval with Reduced Communication and Computation

Zeyu Liu¹, Eran Tromer², Yunhao Wang¹

¹Yale University

²Boston University

November 21, 2024

Abstract

Anonymous message delivery, as in privacy-preserving blockchain and private messaging applications, needs to protect recipient metadata: eavesdroppers should not be able to link messages to their recipients. This raises the question: how can untrusted servers assist in delivering the pertinent messages to each recipient, without learning which messages are addressed to whom?

Recent work constructed Oblivious Message Retrieval (OMR) protocols that outsource the message detection and retrieval in a privacy-preserving way, using homomorphic encryption. Their construction exhibits significant costs in computation per message scanned (~ 0.1 second), as well as in the size of the associated messages (~ 1 kB overhead) and public keys (~ 132 kB).

This work constructs more efficient OMR schemes, by replacing the LWE-based clue encryption of prior works with a Ring-LWE variant, and utilizing the resulting flexibility to improve several components of the scheme. We thus devise, analyze, and benchmark two protocols:

The first protocol focuses on improving the detector runtime, using a new retrieval circuit that can be homomorphically evaluated 13.8x faster than the prior work.

The second protocol focuses on reducing the communication costs, by designing a different homomorphic decryption circuit. While the circuit is less homomorphic-encryption-friendly (than our first construction), it allows the parameter of the Ring-LWE encryption to be set such that both the public key and the message size are greatly reduced. Concretely, the public key size is about 235x smaller than the prior work, and the message size is roughly 1.6x smaller. The runtime of this second construction is ~ 40.0 ms per message, still more than 2.5x faster than prior works.

Contents

1	Introduction	4
1.1	Our Contribution	5
1.2	Related Works	5
2	Technical Overview	7
2.1	Comparisons	8
3	Preliminaries	8
3.1	Hard Problems	9
3.2	PVW Encryption	9
3.3	Fully Homomorphic Encryption	10
3.4	Oblivious Message Retrieval (Definition)	10
4	Revisiting the OMRp2 Construction	12
4.1	Setup	12
4.2	Retrieval	13
4.2.1	Step 1: From Bulletin Board to Pertinency Vector	13
4.2.2	Step 2: Unpack the Pertinency Vector	14
4.2.3	Step 3: Use PV to Construct the Digest	16
4.3	Decoding	17
5	Improved OMR Construction	18
5.1	Reducing the Clue Key Size using RLWE	18
5.2	A New Retrieval Circuit	21
5.2.1	A New Homomorphic Decryption Circuit for Step 1	21
5.2.2	An Efficient PV Unpacking Algorithm for Step 2	22
5.2.3	A General Encoding Procedure for Step 3	24
5.3	Putting Everything Together	28
5.4	Optimizations	29
6	OMR with Further Reduced Key and Clue Sizes	30
6.1	Reducing the Clue Field	30
7	Evaluation	33
7.1	Methodology	33
7.2	Results	33
7.3	Integration Considerations	36
	Acknowledgements	37
	References	38
A	Additional discussion	41
A.1	Boosting Soundness	41
A.2	An Alternative Way to do Index Encoding	41

1 Introduction

Protecting message contents in messaging applications has been extensively studied, with the wide usage of end-to-end encryption today. However, metadata (who sent and received messages, and when) can by itself disclose sensitive information. Therefore, protecting metadata is essential to anonymous message delivery systems like private messaging [45, 13, 29, 8, 4]. The importance is further amplified in privacy-preserving cryptocurrencies [33, 5, 21, 9], where the underlying ledger containing all the messages is permissionless, decentralized, and widely replicated, making it easily accessible to everyone.

Among the various critical pieces of metadata, *recipient* privacy is a particularly challenging problem to efficiently solve. From a recipient’s perspective, a transaction *pertinent* to them could be located anywhere on the ledger in the blockchain applications (or a queue in private messaging systems). Consequently, to find the messages pertinent to them, one simple way is for every recipient to scan the entire ledger. However, the imposed communication and computation costs may be too much of a burden for recipients with limited resources (e.g., wallet apps running on mobile devices). It is desirable to outsource this burden to a server in a privacy-preserving way.

Fuzzy Message Detection (FMD) [4, 41] is the first primitive proposed to address this issue. It adopts a decoy-based approach: the server detects and forwards a set of messages, where the messages pertinent to the recipient are buried among many other additional randomly chosen messages. This is a weak, non-robust security guarantee [43].

Two primitives emerged after FMD enhanced the security guarantee to entirely hide the set of pertinent messages. Private Signaling (PS) [31, 23] focuses on achieving this functionality by leveraging a trusted execution environment (TEE) or two communicating-but-non-colluding servers, while Oblivious Message Retrieval (OMR) [26, 27] realizes it via cryptographic assumptions on a single server. The state-of-the-art PS work [23] provides a very scalable solution; nonetheless, this line of work has a much stronger environmental assumption than OMR. Thus, in this work, we focus on OMR instead of PS.

System Model. OMR works in the following model: in the systems, there are *senders* who send messages to the *recipients* without revealing who the recipients are. Each *message* contains a *payload* and a *clue* generated by the sender using the recipient’s *clue key*. All the messages are placed on a *bulletin board*. When the recipients want to retrieve their messages, they send a *detection key* to an untrusted server, denoted the *detector*. The detector uses the board and the detection key to generate a *digest* and sends it back to the recipients. The recipients decode the digest to obtain all the payloads pertinent to them.

Threat model. We consider an adversary that wishes to learn metadata about which messages are addressed to which user, and about the identity of users that perform message-fetching queries.

The adversary can read all public information (including all board messages and all public keys in the system), and all communication between detectors and the recipients. The adversary may control, or collude with, all the parties in the systems, except for the sender and recipient(s) of the message(s) whose privacy is to be protected. The adversary and its colluding parties may behave maliciously and send malformed messages and keys; but they are computationally bounded (i.e., cannot break the underlying computational assumptions).

Prior OMR schemes. Under these models, OMR [26] (and its extension to group setting [27]) offer a solution based on PVW encryption [39] and BFV homomorphic encryption [10, 15]. However,

there are several major practical concerns about the existing constructions:¹ (1) the detector is slow, taking ~ 109 ms/msg (1-thread, and ~ 51.7 ms/msg 4-thread) per processed message, i.e., approximately \$17.6 per month for Bitcoin-scale application; (2) the clue key size is ~ 132 kB, which is awkward to senders as part of the recipient’s public address/key; (3) the clue itself has a size of ~ 956 bytes, which roughly doubles the total message size, e.g., for typical Zcash transactions.

In this paper, we address these practical issues by presenting two protocols that greatly outperform the prior works, and offer different tradeoffs between computation, clue sizes and key sizes.

1.1 Our Contribution

New constructions. We show two new constructions that provide different trade-offs. Both of the constructions are based on standard lattice hardness assumption (Ring-LWE) and are thus plausibly post-quantum secure.

- For our first construction, we tailor Ring-LWE encryption into a variant that suits our application. Combining this tailored scheme and a newly designed retrieval circuit for the detector (including a new decryption circuit and a new way to encode the digest), we obtain a construction that is both asymptotically and concretely faster in terms of the detector runtime.
- In the second construction, we show an alternative way to parametrize our Ring-LWE variant, together with a new decryption circuit. This alternative construction achieves a much smaller clue and clue key size, and a detector runtime that is still faster than prior work [26].

Implementation and evaluation. We implement our constructions in a C++ library [36] and measure the concrete performance improvement. Salient observations include:

- With the first construction, the detector runtime is about 13.8x faster, and the clue key size is about 60x smaller. The detector runtime is only ~ 7.9 thread-ms/msg scanned (~ 3.7 ms/msg 2-thread), thus only costs $\sim \$0.12$ per million message scanned ($\$1.12$ /month for Bitcoin-scale applications).
- With the second construction, the detector runtime is about 2.7x faster than prior work (~ 40.0 ms/msg 1-thread, ~ 20.2 ms/msg 2-thread), while the clue key size is about 235x smaller. Furthermore, the clue size is about 1.6x smaller. These advantages allow the applications to have much less clue distribution and message size burdens.

We also discuss implications of these improvements on integration with a blockchain-based privacy-preserving cryptocurrency (exemplified by Zcash).

1.2 Related Works

Oblivious Message Retrieval. OMR [26] first proposes a message retrieval primitive with full recipient privacy. Later GOMR [27] extends it to the group setting. Both works rely on a hybrid use of the PVW encryption scheme and BFV leveled homomorphic encryption scheme. We recap

¹Benchmarked using the parameters of [27, Sec 9] and Google Cloud prices (instance type `e2-standard-4`), amortized.

	ClueToPackedPV		PVUnpack		ExpandedPVToDigest		Clue Size	Clue Key Size	Detection Key Size	Digest Size
	# of hom. operations	Depth	# of hom. operations	Depth	# of hom. operations	Depth				
OMRp2 [26, 27]	$O(N\ell t)$	$O(\log(\ell t))$	$O(N \log D)$ or $O(N)$	1 or $\log(D)$	$\tilde{O}(P \cdot N)$	1	$O(n \log(t))$	$\omega(n\ell \log(n) \log(t))$	poly in homomorphic circuit depth	$\tilde{O}(P(\bar{k} + N\epsilon_p))$
PerfOMR1Section 5	$O(N\ell(\log(t) + h))$	$O(\log(\ell th))$					$O(n \log(t))$	$O(n \log(t))$		
PerfOMR2Section 6	$O(N\ell(q \cdot h))$	$O(\log(\ell qh))$	$O(N/v)$	1			$O(n \log(q))$	$O(n \log(q))$		$\tilde{O}(P(\bar{k} + N\epsilon_p)v)$

Table 1: Asymptotic comparison with prior construction. N is the total number of messages. \bar{k} is the upper bound of the number of pertinent messages provided by the recipient during retrieval. ϵ_p is the false positive rate. n, ℓ, q, h are all PVW encryption or sRLWE scheme parameters (see Section 3.2 and Section 5.1). t is the BFV plaintext space. Practically $t \geq qh \gg \log(t) + h$, and D is the BFV ring dimension. P is the size of a payload (which can also be viewed as a constant). v is a tune-able parameter in our construction, which essentially means “gluing” v messages together and treating them as a single message in the later phases during detection (see Sections 4.2.2 and 4.2.3).

how the construction of [26] works in Section 4, and compare our schemes with it asymptotically in Section 2.1 and concretely in Section 7.

Fuzzy Message Detection. FMD [4, 41] mainly focuses on decoy based security. While this primitive has highly efficient constructions, we consider the security notion is relatively weak as analyzed in [43]. Therefore, we do not compare these constructions directly.

Private Signaling. Like OMR, PS [31, 23] provides full security. However, prior works on private signaling have constructions using a Trusted Execution Environment (TEE). TEE is a strong environment assumption since a lot of work shows that the existing TEEs have side-channels that can leak secrets easily [44]. Therefore, while the construction in [23] provides a construction with great scalability (the runtime growth is only poly-logarithmic in the number of messages), we do not directly compare to them as we are assuming a standard environment.

[31] also provides a solution assuming two communicating but non-colluding servers, which is also a very strong environment assumption. Moreover, this construction also scales linearly the number of messages and is concretely slower than [26] and thus our constructions. Therefore, we do not compare with this construction directly either.

PIR. Other related problems are *Private Information Retrieval (PIR)* [12] and its variant *Keyword PIR* [11]. ; and in particular, since OMR recipients retrieve multiple messages, the most related primitive is the variant called *multi-query (keyword) PIR* or *batch (keyword) PIR*. Our setting differs in that recipients do not know the indices or labels of messages pertinent to them; rather, the clues are randomized and require nontrivial computation (rather than simple comparison) to detect.

Private Stream Search. In Private Stream Search (PSS) [37, 14, 7, 16], a client can search a keyword over a database of documents and retrieve the ones with such a keyword without revealing the keyword to the server. As for Keyword PIR, this does not directly yield OMR. In [26], the authors use similar techniques as in PSS works, for the index encoding. While our scheme builds upon [26], our index encoding is different, requiring additional techniques to handle.

See [26, 27] for further discussion.

2 Technical Overview

We follow the OMR framework introduced in [26] to build our improved constructions, and the requisite background is systematically recalled in Section 3 and Section 4; then the improvements are presented in detail in Section 5 and Section 6. For those readers already familiar with the approach of [26] and the encryption schemes it employs (PVW [39] and BFV [10, 15]), the following succinctly summarizes our approach to improvements.

Setup with tailored RLWE Encryption. Whereas [26] had clues which are PVW encryptions (based on LWE hardness), we instead use an encryption scheme based on Ring-LWE hardness. Our encryption scheme, `sRLWE`, is a variant of RLWE [30] but with sparse key, smaller decryption range, and smaller plaintext space. The encryption public key `sRLWE.pk` (included in the clue key) is much smaller than with PVW.

The sender generates $\text{sRLWE.Enc}(\text{sRLWE.pk}, 0) \in \mathbb{Z}_t^{n+1}$ as the clue (for some security parameter n , ciphertext modulus q)². To perform a retrieval, the recipient uses the homomorphic encryption scheme BFV to compute $\text{ct}_{\text{sk}} \leftarrow \text{BFV.Enc}(\text{BFV.pk}, \text{sRLWE.sk})$ and send $(\text{BFV.pk}, \text{ct}_{\text{sk}})$ as the detection key to the detector.

A more efficient homomorphic decryption circuit. Given a detection key, the first step performed by the detector is to homomorphically decrypt each `sRLWE` ciphertext over \mathbb{Z}_t . Prior work relies on a degree- $(t-1)$ polynomial as in [26], which requires $t-1$ homomorphic operations. By exploiting the fact that `sRLWE` relies on sparse secret-key RLWE (using secrets with fixed hamming weight h), which implies that the `sRLWE` ciphertexts have noise $O(h)$, we design a more efficient decryption circuit that only takes $O(h + \log(t))$ operations. Since this circuit is evaluated using BFV, D (BFV ring dimension) clues are homomorphically decrypted simultaneously. For N clues, this process results in $d = \lceil N/D \rceil$ BFV ciphertexts, each of which encrypts a binary vector of size D (in its D slots) representing whether D corresponding messages are pertinent. We call this step `ClueToPackedPV`.

A new way to expand the BFV ciphertexts. The next step for the detector is to homomorphically expand these ciphertexts. Instead of one ciphertext encrypting D bits, each bit represents whether a message is pertinent, the detector needs D ciphertexts each encrypting a single bit repeated D times (for more efficient digest encoding). To accomplish this goal, we first homomorphically decode the ciphertext via the `SlotToCoeff` procedure in [28]. Then, we perform `OExpand` introduced in [3] on the decoded ciphertexts to obtain the targeted result. This new expansion way requires only $O(D)$ homomorphic operations for each ciphertext, compared to $O(D \log(D))$ operations in prior work [26]. We call this step `PVUnpack`.

Bundling \mathbf{v} messages. With this new way of expansion, it still takes $O(N)$ homomorphic operations for N messages. To further reduce the cost, a natural way is to bundle \mathbf{v} messages to a single message. This can be done by adding up \mathbf{v} ciphertexts obtained from `ClueToPackedPV` before expanding the ciphertexts.

A new encoding scheme. Despite the improved efficiency of the bundling technique, it introduces extra complexity. The major issue is that the encoding scheme in [26] does not work anymore, given that the ciphertexts output from `PVUnpack` now encrypts non-binary values (since we add \mathbf{v} binary values together). To resolve this, we design a new encoding scheme for index encoding:

²In the actual construction, we encrypt 0^ℓ for some $\ell \geq 1$, to reduce false positive rate. We set $\ell = 1$ here for simplicity.

we first expand each bit of the indices into $\log(v + 1)$ bits; then we use these expanded indices to encode and allow the recipient to decode all the pertinent indices. We call this last encoding step `ExpandedPVToDigest`

Putting all these three steps `ClueToPackedPV`, `PVUnpack`, `ExpandedPVToDigest` together, we obtain our first construction `PerfOMR1`, which is both asymptotically and concretely faster than the prior work `OMRp2` in [26].

An alternative way to use sRLWE. Another way to use sRLWE is that instead of having its ciphertext modulus be t (the same as the BFV plaintext modulus, which relatively large for practicality), we can set sRLWE modulus to $q \ll t$. As our sRLWE relies on sparse keys (keys with hamming weight h), we set $qh < t/2$. This guarantees that decrypting the sRLWE ciphertext over \mathbb{Z}_t is the same as over \mathbb{Z} (no wrap-arounds). Therefore, we can instead design a polynomial with $O(qh)$ degree to perform the homomorphic decryption. While this makes the runtime worse, the clue key and clue of size $O(n \log(q))$ can be greatly reduced as q now is smaller.

2.1 Comparisons

In Table 1, we compare the asymptotic behavior of our constructions, in terms of the cost metrics, with the prior construction in [26, 27]. As mentioned in above, our work mainly focuses on the improvement of the detector construction, which is composed of three main steps: `ClueToPackedPV`, `PVUnpack`, `ExpandedPVToDigest`.

For our first construction, `PerfOMR1`, the detector runtime is strictly faster than the prior works by having much fewer homomorphic operations: in the step `ClueToPackedPV`, h is the hamming weight of the secret key which is normally viewed as $O(1)$ in terms of security parameter (e.g., [18]), and we thus have $\log(t) + h = o(t)$; in the step `PVUnpack`, we have $v \geq 1$, thus reducing the number of homomorphic operations by a factor of v . Besides, the clue key size is smaller by reducing $\omega(n \log(n))$ to $O(n)$.

For our second construction `PerfOMR2`, we set $q \cdot h \leq t$. Therefore, the runtime is comparable with the prior work with slightly fewer homomorphic operations in the step `ClueToPackedPV`. The gain is that the clue size and the clue key size are both smaller since they now depend on $q < t$.

Note that the digest size of both of our constructions is parametrized by v . Concretely, the digest size is exactly the same as the prior work when $v = 1$. Since v only affects the runtime of `PVUnpack` step, when `PVUnpack` is the runtime bottleneck, we set $v > 1$ (e.g., for `PerfOMR1`, we set $v = 8$ to reach the optimal runtime); otherwise, we set $v = 1$ (e.g., for `PerfOMR2`).

See Section 7 for evaluation of concrete performance.

3 Preliminaries

Notation. Let $[n]$ denote the set $\{1, \dots, n\}$. For a vector \mathbf{v} , $\mathbf{v}[i]$ indicates the i -th element of this vector. For a matrix A , $A[i][j]$ indicates the cell at the i -th row and j -th column. Let $\mathcal{R} = \mathbb{Z}[X]/(X^N + 1)$ denote the $2N$ -th cyclotomic ring where N is a power-of-two, and $\mathcal{R}_Q = \mathcal{R}/Q\mathcal{R}$ for some $Q \in \mathbb{Z}$. For matrices $A, B \in \mathbb{Z}_t^{n \times m}$, let \circ denote the Hadamard product $C \leftarrow A \circ B$ satisfying $C[i][j] = A[i][j] \cdot B[i][j], \forall i \in [n], j \in [m]$. Drawing x uniformly at random from a set S is denoted $x \xleftarrow{\$} S$.

3.1 Hard Problems

Definition 3.1 (Decisional learning with error problem). Let n, q, \mathcal{D}, χ be parameters dependent on λ . The learning with error (LWE) problem states the following: for $a \xleftarrow{\$} \mathbb{Z}_q^n$, it holds that $(a, \langle a, s \rangle + e) \approx_c (a, b)$, where $s \leftarrow \mathcal{D}, e \leftarrow \chi$ and $b \xleftarrow{\$} \mathbb{Z}_q$.

Let $\text{LWE}_{n,q,\mathcal{D},\chi}$ denote the LWE assumption parameterized by n, q, \mathcal{D}, χ .

Definition 3.2 (Decisional ring learning with error problem). [30, 40] Let N, Q, \mathcal{D}, χ be parameters dependent on λ and N being a power of two. Let $\mathcal{R} = \mathbb{Z}[X]/(X^N + 1)$. The ring learning with error (RLWE) problem $\text{RLWE}_{N,Q,\mathcal{D},\chi}$ is the following: distinguish $(a, a \cdot s + e)$ and (a, b) (with noticeable advantage), where $a \xleftarrow{\$} \mathcal{R}_Q, s \leftarrow \mathcal{D}, e \leftarrow \chi$ and $b \xleftarrow{\$} \mathcal{R}_Q$.

3.2 PVW Encryption

We adapt the PVW encryption from [39] and modify it according to [26].

- $\text{pp} = (n, \ell, w, q, \sigma, r) \leftarrow \text{PVW.GenParam}(1^\lambda, \ell, q, \sigma, \epsilon_n)$: Choose a secret key dimension n , and $w = \omega(n \log(q))$ by ciphertext modulus q , plaintext size ℓ , and standard deviation σ for Gaussian distribution for ciphertext noise generation, as in [39]. Additionally, choose the noise bound r such that $\Pr[\text{PVW.Dec}(\text{sk}, \text{PVW.Enc}(\text{pk}, \vec{m})) = \vec{m}] \geq 1 - \epsilon_n - \text{negl}(\lambda)$.
- $(\text{sk}, \text{pk}) \leftarrow \text{PVW.KeyGen}(\text{pp})$: Draw a secret key $\text{sk} \xleftarrow{\$} \mathbb{Z}_q^{n \times \ell}$. Sample $A \xleftarrow{\$} \mathbb{Z}_q^{n \times w}$ and a noise matrix $X \in \mathbb{Z}_q^{\ell \times w}$ from the Gaussian distribution χ_σ , and compute $\text{pk} = (A, P = \text{sk}^T A + X)$.
- $\text{ct} = (\vec{a}, \vec{b}) \leftarrow \text{PVW.Enc}(\text{pp}, \text{pk} = (A, P), \vec{m})$: To encrypt a vector $\vec{m} \in \mathbb{Z}_2^\ell$, define the vector $\vec{t} = \frac{q}{2} \cdot \vec{m} \in \mathbb{Z}_q^\ell$, and draw $\vec{e} \xleftarrow{\$} \{0, 1\}^w \in \mathbb{Z}_2^w$. The ciphertext is the pair $(\vec{a}, \vec{b}) = (A\vec{e}, P\vec{e} + \vec{t}) \in \mathbb{Z}_q^n \times \mathbb{Z}_q^\ell$.
- $\vec{m} \leftarrow \text{PVW.Dec}(\text{pp}, \text{sk}, \text{ct} = (\vec{a}, \vec{b}))$: $\vec{d} = \vec{b} - \text{sk}^T \vec{a} \in \mathbb{Z}_q^\ell$, let $\vec{m} \in \mathbb{Z}_2^\ell$, and $\vec{m}[i] = 1$ iff $\vec{d}[i] + r/2 > r$ for all $i \in [\ell]$.

The scheme satisfies CPA security and tailored correctness: correct with probability $1 - \epsilon_n$ for some $0 < \epsilon_n < 1$. The public key size is $\omega(\ell n \log^2(q))$ (the size of P in bits, as A can be represented by a random seed).

PVW also has the following properties:

1. (Key privacy) Two ciphertexts encrypted under two different public keys are computationally indistinguishable. Formally speaking, for any PPT adversary $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$, for any $\lambda > 0$, $\ell = \text{poly}(\lambda), \sigma > 0, 1 > \epsilon_n > 0, q = \text{poly}(\lambda, \sigma, \epsilon_n)$, let $\text{pp}_{\text{PVW}} \leftarrow \text{PVW.GenParam}(1^\lambda, \ell, q, \sigma, \epsilon_n)$, $(\text{sk}, \text{pk}) \leftarrow \text{PVW.KeyGen}(\text{pp}_{\text{PVW}})$, $(\text{sk}', \text{pk}') \leftarrow \text{PVW.KeyGen}(\text{pp}_{\text{PVW}})$. The adversary then chooses a message and remembers its state $(m, \text{st}) \leftarrow \mathcal{A}_1(\text{pp}_{\text{PVW}}, \text{pk}, \text{pk}')$. Let $\text{ct} \leftarrow \text{PVW.Enc}(\text{pp}_{\text{PVW}}, \text{pk}, m)$, $\text{ct}' \leftarrow \text{PVW.Enc}(\text{pp}_{\text{PVW}}, \text{pk}', m)$, it holds that: $|\Pr[\mathcal{A}_2(\text{st}, \text{ct}) = 1] - \Pr[\mathcal{A}_2(\text{st}, \text{ct}') = 1]| \leq \text{negl}(\lambda)$.
2. (Zero-plaintext wrong-key decryption) Given the wrong key, a PVW ciphertext is decrypted into a zero plaintext with probability $\leq (q-r)^{-\ell} + \text{negl}(\lambda)$. Formally speaking: let $\text{pp}_{\text{PVW}}, (\text{sk}, \text{pk}), (\text{sk}', \text{pk}')$ be generated as above, and $\text{ct} = \text{PVW.Enc}(\text{pk}, 0^\ell)$, it holds that $\Pr[\text{PVW.Dec}(\text{sk}', \text{ct}) = 0^\ell] \leq (q-r)^{-\ell} + \text{negl}(\lambda)$.

We refer readers to the formal proof of these two properties in [26].

3.3 Fully Homomorphic Encryption

Fully Homomorphic Encryption (FHE), introduced by Rivest et al. [42] and first constructed by Gentry [19], enables evaluation of a circuit on encrypted data, such that the result is the encryption of the corresponding output.

BFV FHE scheme. We use the Brakerski/Fan-Vercauteran (BFV) homomorphic encryption scheme [10, 15] in all constructions.

BFV scheme consists of the following PPT algorithms: $\text{GenParam}(1^\lambda), \text{KeyGen}(\text{pp}_{\text{BFV}}), \text{Enc}(\text{pp}_{\text{BFV}}, \text{pk}, m), \text{Dec}(\text{pp}_{\text{BFV}}, \text{sk}, c)$ as normal PKE schemes. BFV is unconditionally correct and sound. Under the *RLWE* hardness assumption, it also fulfills the standard definitions of semantic security (IND-CPA) for FHE schemes.

Given a polynomial from the cyclotomic ring $R_t = \mathbb{Z}_t[X]/(X^D + 1)$ (where D is a power-of-two, $t \equiv 1 \pmod{2D}$), the BFV scheme encrypts it into a ciphertext consisting of two polynomials, each of which in a larger cyclotomic ring $R_Q = \mathbb{Z}_Q[X]/(X^D + 1)$ for some $Q > t$. Here, t , Q , and D are called the plaintext modulus, the ciphertext modulus, and the ring dimension, respectively.

Plaintext encoding. In practice, instead of having a polynomial in $\mathcal{R}_t = \mathbb{Z}_t[X]/(X^D + 1)$ directly as input, applications usually hold a vector of messages $\vec{m} = (m_1, \dots, m_D) \in \mathbb{Z}_t^D$. Thus, to encrypt such input messages, BFV first encodes the messages into another polynomial $y(X) = \sum_{i \in [D]} y_i X^{i-1}$ such that $m_j = y(\zeta_j)$, $\zeta_j := \zeta^{3^j} \pmod{t}$, and ζ is the $2N$ -th primitive root of unity of t . Such encoding can be done using Inverse Number Theoretic Transform (INTT), which is a linear transformation represented as matrix multiplication. We say that a BFV ciphertext has D slots, each of which is a \mathbb{Z}_t element.

For simplicity, we assume BFV.Enc takes a vector of form \mathbb{Z}_t^D as an input, and BFV.Dec outputs a vector of form \mathbb{Z}_t^D , and will handle encode and decode implicitly. We use BFV.PartialDec to represent decryption without decoding. In other words, for a ciphertext ct , the output of $\text{BFV.PartialDec}(\text{sk}, \text{ct}) \in \mathcal{R}_t$ is the encoding of $\text{BFV.Dec}(\text{sk}, \text{ct}) \in \mathbb{Z}_t^D$.

Operations. BFV supports the following operations.

- (Additions) For any two BFV ciphertexts ct_1, ct_2 , and $\text{ct} \leftarrow \text{ct}_1 + \text{ct}_2$, it holds that $\text{BFV.Dec}(\text{ct}) = \text{BFV.Dec}(\text{ct}_1) + \text{BFV.Dec}(\text{ct}_2)$ (element-wise).
- (Multiplication) For any two BFV ciphertexts ct_1, ct_2 , and $\text{ct} \leftarrow \text{ct}_1 \times \text{ct}_2$, it holds that $\text{BFV.Dec}(\text{ct}) = \text{BFV.Dec}(\text{ct}_1) \times \text{BFV.Dec}(\text{ct}_2)$ (element-wise).
- (Rotation) For any BFV ciphertexts ct , and $\text{ct}' \leftarrow \text{BFV.Rotate}(\text{ct}, k)$ for some $k \in [D]$, let $\text{BFV.Dec}(\text{sk}, \text{ct})[i] = \text{BFV.Dec}(\text{sk}, \text{ct}')[i + k \pmod{D}]$.
- (Substitution) For any BFV ciphertexts ct , and $\text{ct}' \leftarrow \text{BFV.Substitute}(\text{ct}, k)$ for some odd number k , let $y(X) = \text{BFV.PartialDec}(\text{sk}, \text{ct})$ and $y'(X) = \text{BFV.PartialDec}(\text{sk}, \text{ct}')$, it holds that $y'(X) = y(X^k) \in \mathcal{R}_t$.

3.4 Oblivious Message Retrieval (Definition)

We adapt the definition of OMR from [26] by introducing a new parameter ν to relax the soundness and compactness as follows. At a high level, the scheme is allowed to include impertinent payloads,

as long as the final output is still bounded by ν . For example, the scheme can bundle ν payloads together. If a bundle contains a pertinent message, the scheme can return all ν payloads in the bundle to the recipient. For example, if the first message is pertinent, the final output of the OMR scheme might contain messages $[1, \nu]$ to the recipient, provided that messages $[1, \nu]$ are in the same bundle. We discuss why this relaxation is reasonable in more detail in Remark 3.5.

Definition 3.3 (Oblivious Message Retrieval(OMR)). An Oblivious Message Retrieval scheme has the following PPT algorithms:

- $\text{pp} \leftarrow \text{GenParam}(1^\lambda, \epsilon_p, \epsilon_n)$: takes a security parameter λ , a false positive rate ϵ_p , a false negative rate ϵ_n , and outputs a public parameter pp .
- $(\text{sk}, \text{pk} = (\text{pk}_{\text{clue}}, \text{pk}_{\text{detect}})) \leftarrow \text{KeyGen}(\text{pp})$: takes the public parameter pp ; outputs a secret key sk and a public key pk consisting of a clue key pk_{clue} and a detection key $\text{pk}_{\text{detect}}$.
- $c \leftarrow \text{GenClue}(\text{pp}, \text{pk}_{\text{clue}}, x)$: takes the public parameter pp , a clue key pk_{clue} , and a payload $x \in \mathcal{P}$ where $\mathcal{P} := \{0,1\}^P$ for some $P > 0$; outputs a clue $c \in \mathcal{C}$.
- $M \leftarrow \text{Retrieve}(\text{pp}, \text{BB}, \text{pk}_{\text{detect}}, \bar{k})$: takes the public parameter pp , a bulletin board $\text{BB} = \{(x_1, c_1), \dots, (x_N, c_N)\}$ for size N , a detection key $\text{pk}_{\text{detect}}$, and an upper bound \bar{k} on the number of pertinent messages addressed to that recipient; outputs a digest M .
- $\text{PL} \leftarrow \text{Decode}(\text{pp}, M, \text{sk})$: takes the public parameter pp , the digest M and the corresponding secret key sk ; outputs either a decoded payload list $\text{PL} \subset \mathcal{P}^k$ or an overflow indication $\text{PL} = \text{overflow}$.

To define soundness and completeness, we first define the notion of board generation:

Definition 3.4 (Board Generation). Given pp , and the size of bulletin board N : arbitrarily choose the number of recipients $1 \leq p \leq N$, and a partition of set $[N]$ into p subsets S_1, \dots, S_p representing the indices of messages addressed to each party. Also arbitrarily choose unique payloads (x_1, \dots, x_N) . For each recipient $i \in [p]$: generate keys $(\text{sk}_i, \text{pk}_i = (\text{pk}_{\text{clue}_i}, \text{pk}_{\text{detect}_i})) \leftarrow \text{KeyGen}(\text{pp})$, and for each $j \in S_i$, generate $c_j \leftarrow \text{GenClue}(\text{pk}_{\text{clue}_i}, x_j)$. Then, output the board $\text{BB} = \{(x_1, c_1), \dots, (x_N, c_N)\}$, the set S_1 , and $(\text{sk}_1, \text{pk}_1 = (\text{pk}_{\text{clue}_1}, \text{pk}_{\text{detect}_1}))$.³

The scheme must satisfy the following properties:

- (Completeness) Let $\text{pp} \leftarrow \text{GenParam}(1^\lambda, \epsilon_p, \epsilon_n, \nu)$. Set any $N = \text{poly}(\lambda)$, and $0 < \bar{k} \leq N$. Let a board BB , a set S of pertinent messages, and a key pair $(\text{sk}, \text{pk} = (\text{pk}_{\text{clue}}, \text{pk}_{\text{detect}}))$ be generated as in Definition 3.4 for any choice of p , partition and payloads therein. Let $M \leftarrow \text{Retrieve}(\text{BB}, \text{pk}_{\text{detect}}, \bar{k})$ and $\text{PL} \leftarrow \text{Decode}(M, \text{sk})$. Let $k = |S|$ (the number of pertinent messages in S). Then either $k > \bar{k}$ and $\text{PL} = \text{overflow}$, or:

$$\Pr[x_j \in \text{PL} \mid j \in S] \geq (1 - \epsilon_n - \text{negl}(\lambda)) \quad \text{for all } j \in [N] \quad .^4$$

- (ν -Soundness) For the same quantifiers as in Completeness:

$$|\text{PL}| = \tilde{O}(\nu \cdot P \cdot (\bar{k} + \epsilon_p N))$$

³That is, S_1 is the indices of messages pertinent to the recipient whose keys are sk_1, pk_1 , which wlog is the first recipient.

- (Computational privacy) For any PPT adversary $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$: let $\text{pp} \leftarrow \text{GenParam}(\epsilon_p, \epsilon_n)$, $(\text{sk}, \text{pk} = (\text{pk}_{\text{clue}}, \text{pk}_{\text{detect}})) \leftarrow \text{KeyGen}(\text{pp})$ and $(\text{sk}', \text{pk}' = (\text{pk}'_{\text{clue}}, \text{pk}'_{\text{detect}})) \leftarrow \text{KeyGen}(\text{pp})$. Let the adversary choose a payload x and remember its state: $(x, \text{st}) \leftarrow \mathcal{A}_1(\text{pp}, \text{pk}, \text{pk}')$. Let $c \leftarrow \text{GenClue}(\text{pk}_{\text{clue}}, x)$ and $c' \leftarrow \text{GenClue}(\text{pk}'_{\text{clue}}, x)$. Then:

$$|\Pr[\mathcal{A}_2(\text{st}, c) = 1] - \Pr[\mathcal{A}_2(\text{st}, c') = 1]| \leq \text{negl}(\lambda) .$$

An OMR scheme is ν -compact if it moreover satisfies the following:

- (ν -Compactness) An OMR scheme is ν -compact if for $\text{pp} \leftarrow \text{GenParam}(1^\lambda, \epsilon_p, \epsilon_n)$, $(\text{sk}, \text{pk} = (\text{pk}_{\text{clue}}, \text{pk}_{\text{detect}})) \leftarrow \text{OMR.KeyGen}(\text{pp})$, for any board $\text{BB} = \{(x_1, c_1), \dots, (x_N, c_N)\}$, letting $M \leftarrow \text{Retrieve}(\text{BB}, \text{pk}_{\text{detect}}, \bar{k})$, it always holds that:

$$|M| = \text{poly}(\lambda, \log N) \cdot \log \epsilon_p^{-1} \cdot \tilde{O}(\bar{k} + \epsilon_p N) \cdot \nu .$$

In the compactness definition, $\tilde{O}(\bar{k} + \epsilon_p N)$ (where $\tilde{O}(x) = x \cdot \text{polylog}(x)$) accounts for the number of messages detected as pertinent, including false positives; and the remaining factors account for the cost of representing each such message, taking the payload size as constant.

Remark 3.5. Note that we have relaxed the “soundness” and “compactness” properties in [26, Def 4.3] special case to “ ν -soundness” and “ ν -compactness”, allowing the digest and the decoded payloads to include more than just the pertinent messages, by a factor of ν . The scheme is thus allowed to bundle $O(\nu)$ messages as a single one and process them together (where each bundle may include pertinent messages and impertinent ones simultaneously).

In most applications of OMR, the recipients are able to find the single intended data payload from a bundle (e.g., the payloads are encrypted and using the wrong key to decrypt is detectable as a decryption failure). Therefore, in many cases, it is not an issue. Section 5.2.3 and Appendix A.1 also shows a general way to guarantee full soundness with a small cost.

4 Revisiting the OMRp2 Construction

We first revisit and summarize the construction of OMR, OMRp2 in [26, Alg 8], which is the basis for improvements in later sections. Here, we abstract out each step of OMRp2 and provides modular analysis to each step to make the entire framework easier to understand.

4.1 Setup

OMRp2 mainly relies on the PVW encryption (see Section 3.2) for clues and BFV homomorphic encryption scheme (see Section 3.3) for retrieval.

GenParam. Public parameter generation is straightforward. It outputs a public parameter pp including the PVW parameters $\text{pp}_{\text{PVW}} = (n, w, q, \ell, \mathcal{D}, \sigma)$, the BFV parameters $\text{pp}_{\text{BFV}} = (D, t, \dots)$ ⁵, false positive rate ϵ_p and false negative rate ϵ_n .

⁵Technically, we should also set Q , the ciphertext modulus, which is used to guarantee that there is enough noise budget to evaluate the entire circuit. However, it is not used in constructions explicitly, we leave it implicit for better readability.

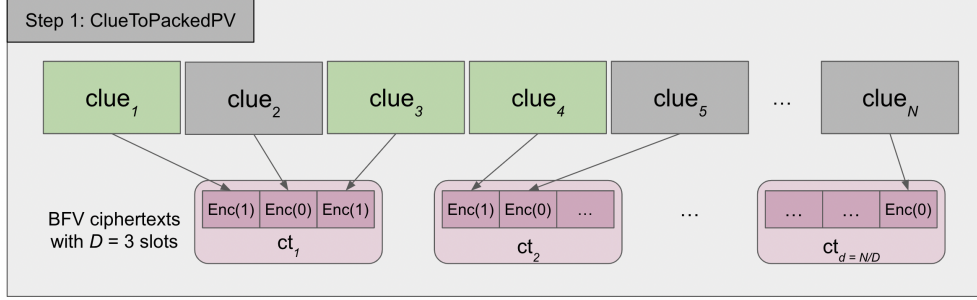


Figure 1: Visualization of Step 1 ClueToPackedPV. The green clues are for the pertinent messages and the gray ones are for the impertinent messages. Each BFV ciphertext in pale pink has $D = 3$ slots and each slot in dark pink encrypts the pertinency of a single message. For N messages, the output has $d = N/D = N/3$ ciphertexts.

KeyGen. The recipient first generates a PVW key pair $(\text{sk}_{\text{pvw}}, \text{pk}_{\text{pvw}})$ and a BFV key pair $(\text{sk}_{\text{BFV}}, \text{pk}_{\text{BFV}})$. pk_{pvw} will be the clue key. The recipient then generates $\text{ct}_{\text{sk}} \leftarrow \text{BFV.Enc}(\text{pp}_{\text{BFV}}, \text{pk}_{\text{BFV}}, \text{sk}_{\text{pvw}})$, which is the encrypted sk_{pvw} under BFV public key pk_{BFV} . The tuple $(\text{ct}_{\text{sk}}, \text{pk}_{\text{BFV}})$ serves as the detection key.

GenClue. After fetching the recipient's pk_{clue} , the sender computes $c \leftarrow \text{PVW.Enc}(\text{pp}_{\text{PVW}}, 1^\ell)$. If a clue is decrypted to 1^ℓ , it indicates that the message is pertinent. Otherwise, there is at least one zero among the ℓ bits, and the message is impertinent. Based on the wrong-key decryption property, if a message is impertinent, the decrypted message will not be 1^ℓ with high probability.

4.2 Retrieval

Recall that the heavy computation work of retrieval is leveraged to a detector. To retrieve the pertinent messages for the recipients, the detector invokes `OMRp2.Retrieve`, which is composed of three main steps. We first define those steps and describe how `OMRp2` in [26] realizes them. Looking ahead, our new construction rewrites these three with better efficiency, both asymptotically and concretely.

4.2.1 Step 1: From Bulletin Board to Pertinency Vector

The first step takes the detection key and all the clues published on the bulletin board, and outputs a vector of BFV ciphertexts, each slot of which indicates whether a single message is pertinent (we call *pertinency vector*, PV). We visualize it in Fig. 1. The interface is defined as follows:

- $(\text{ct}_1, \dots, \text{ct}_d) \leftarrow \text{ClueToPackedPV}(\text{pp}, \text{pk}_{\text{detect}}, \text{BB})$: takes public parameter pp , a detection key $\text{pk}_{\text{detect}}$, and a bulletin board BB of size N ; outputs a vector of BFV ciphertexts $(\text{ct}_1, \dots, \text{ct}_d)$ where $d = \lceil N/D \rceil$.

Recall that each BFV ciphertext contains D slots (i.e., encrypting a vector of \mathbb{Z}_t^D for t being the plaintext modulus), where D is the ring dimension. If the i -th message is pertinent, the i -th slots should be 1; and 0 otherwise. Thus, there are in total of $\lceil N/D \rceil$ ciphertexts with $\geq N$ slots to encrypt the pertinency of each message. Correctness is defined as follows:

Definition 4.1 (Correctness of ClueToPackedPV). Let $\text{pp} \leftarrow \text{GenParam}(1^\lambda, \epsilon_p, \epsilon_n)$. For any $N = \text{poly}(\lambda)$, and $0 < \bar{k} \leq N$, let a board BB , a set S of pertinent messages, and a key pair $(\text{sk}, \text{pk} = (\text{pk}_{\text{clue}}, \text{pk}_{\text{detect}}))$ be generated as in Definition 3.4 for any choice of p , partition and payloads therein, let $(\text{ct}_1, \dots, \text{ct}_d) \leftarrow \text{ClueToPackedPV}(\text{pp}, \text{pk}_{\text{detect}}, \text{BB})$, it holds that:

$$\Pr [\text{BFV.Dec}(\text{sk}, \text{ct}_j)[i] = 1 \mid j \cdot D + i \in S] \geq (1 - \epsilon_n - \text{negl}(\lambda)) \quad \text{for all } i \in [D], j \in [d] .$$

and:

$$\Pr [\text{BFV.Dec}(\text{sk}, \text{ct}_j)[i] = 1 \mid j \cdot D + i \notin S] \leq (\epsilon_p + \text{negl}(\lambda)) \quad \text{for all } i \in [D], j \in [d] .$$

ClueToPackedPV Implementation. In the construction of OMRp2, the detector uses ct_{sk} to homomorphically decrypt each clue (where ct_{sk} is the encryption of the PVW secret key and the clue is a PVW ciphertext). If a message is indeed pertinent, the homomorphic decryption yields 1^ℓ except with $\epsilon_n + \text{negl}(\lambda)$ probability. Otherwise, the result would be 1^ℓ with probability $\leq \epsilon_p + \text{negl}(\lambda)$. Lastly, the detector multiplies all the ℓ bits, and gets 1 if and only if the message is pertinent.

This homomorphic decryption circuit is evaluated under BFV, and D messages are processed simultaneously by taking advantage of the SIMD property of BFV. As mentioned before, the ciphertext ct after the homomorphic decryption has D slots, where the i -th slot encrypts 1 if and only if the i -th message is pertinent (except with some small bounded probability), for $i \in [D]$. This decryption process is repeated $d = \lceil N/D \rceil$ times to obtain ciphertexts for all the N messages.

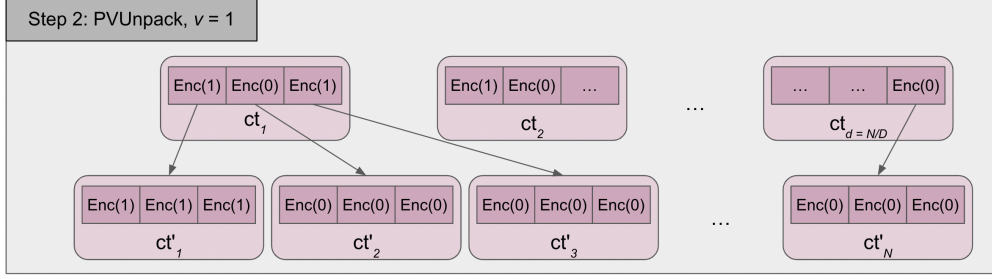
Essentially, let $\text{sk}[i] \in \mathbb{Z}_q^n$ denote the i -th column of $\text{sk} \in \mathbb{Z}_q^{n \times \ell}$, the PVW decryption circuit checks whether $|b[i] - \langle \vec{a}, \text{sk}[i] \rangle| \leq r$ for $i \in [\ell]$, for each clue of form $(\vec{a}, \vec{b}) \in \mathbb{Z}_q^n \times \mathbb{Z}_q^\ell$. The evaluation of $b[i] - \langle \vec{a}, \text{sk}[i] \rangle$ is easy. The hard part is the range check, as BFV operates over a finite field and only supports additions and multiplications. Fortunately, one important observation in [26] is that any function over \mathbb{Z}_t can be represented by a polynomial with degree- $(t-1)$. Therefore, the detector interpolates a degree- $(t-1)$ function to check the range and completes the step ClueToPackedPV.

4.2.2 Step 2: Unpack the Pertinency Vector

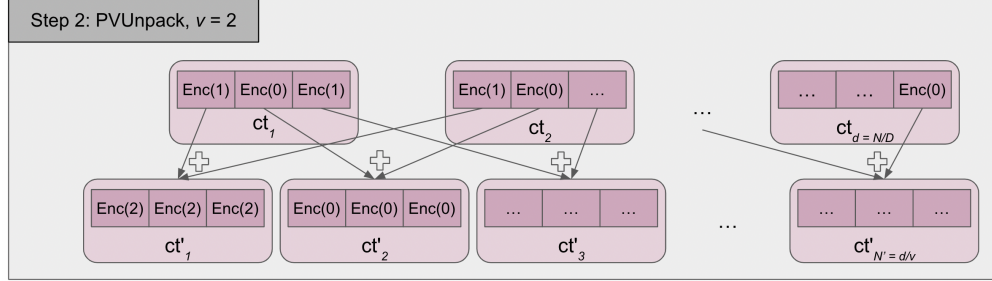
After obtaining the pertinency vector, to prepare for the third step, the framework in [26] unpacks those N slots (of the $\lceil N/D \rceil$ ciphertexts from last step) into N ciphertexts. If the i -th slot is 1, the i -th ciphertext after unpacking encrypts 1 in all of its D slots, and 0 in all D slots otherwise.

We further generalize this step to take a parameter called *bundle size* v to output N/v ciphertexts instead of N ciphertexts. For $v = 1$, this procedure simply unpacks N slots into N separate ciphertexts. Looking ahead, in section Section 5.2.2, we bundle $v > 1$ messages into a single one for better efficiency: the i -th ciphertext after unpacking encrypts the number of slots encrypting 1's among all the v slots corresponding to the bundled messages. For simplicity, we require v to divide d . We visualize step 2 and the difference between $v = 1$ and $v > 1$ in Fig. 2. The interface is as follows:

- $(\text{ct}'_1, \dots, \text{ct}'_{N'}) \leftarrow \text{PVUnpack}(\text{pp}, \text{pk}_{\text{detect}}, (\text{ct}_1, \dots, \text{ct}_d), v)$: takes public parameter pp , a detect key $\text{pk}_{\text{detect}}$, a vector of BFV ciphertexts of length d , the bundle size v (requiring v to divide d for simplicity); outputs a vector of BFV ciphertexts of size $N' = d \cdot D/v$ for D being the underlying BFV ring dimension.



(a) Step 2 with $v = 1$ (i.e., for the original OMR construction in Section 4.2.2)



(b) Step 2 with $v = 2$ (i.e., for our new construction in Section 5.2.2)

Figure 2: Visualization of step 2 PVUnpack. When $v = 1$, each slot is expanded into a single BFV ciphertext, resulting in N ciphertexts. When $v = 2$, two slots are added up and expanded into a single BFV ciphertext, resulting in N/v ciphertexts

The correctness is as follows:

Definition 4.2 (Correctness of PVUnpack). Let $pp, sk, pk = (pk_{\text{clue}}, pk_{\text{detect}})$ generated as in Definition 4.1, for any vector of ciphertexts (ct_1, \dots, ct_d) , let $(ct'_1, \dots, ct'_{N'}) \leftarrow \text{PVUnpack}(pp, pk_{\text{detect}}, (ct_1, \dots, ct_d), v)$, it holds that:

$$\Pr[\text{BFV.Dec}(sk, ct'_{j \cdot D + i}) = \left(\sum_{w=0}^{v-1} \text{BFV.Dec}(sk, ct_{j \cdot v + w})[i] \right)^D] \geq 1 - \text{negl}(\lambda) \quad \text{for all } i \in [D], j \in [0, \frac{d}{v} - 1].^6$$

PVUnpack Implementation (with $v = 1$). We explain detailed construction in [26] by giving a simplified example of unpacking a single ciphertext, i.e., given a BFV ciphertext ct encrypting (b_1, \dots, b_D) where b_i is the i -th encrypted bit, we eventually want D new ciphertexts where the i -th ciphertext encrypts $\vec{b}_i := (b_i, \dots, b_i)$, which is a vector of D b_i 's in all slots.

[26] first multiply ct with $(0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{Z}_t^D$ where the i -th slot is 1 and other slots are 0, obtaining a ciphertext ct' encrypting $(0, \dots, 0, b_i, 0, \dots, 0)$, where b_i is the i -th element encrypted in ct . To fill all the D slots with the same value b_i , [26] uses the rotate-and-add method: for $j \in [\log(D)]$, it computes $ct' \leftarrow ct' + \text{BFV.Rotate}(ct', 2^{j-1})$. The final ciphertext ct' thus encrypts \vec{b}_i as desired. This process is simply repeated D times for each input ciphertext.⁷ This step requires

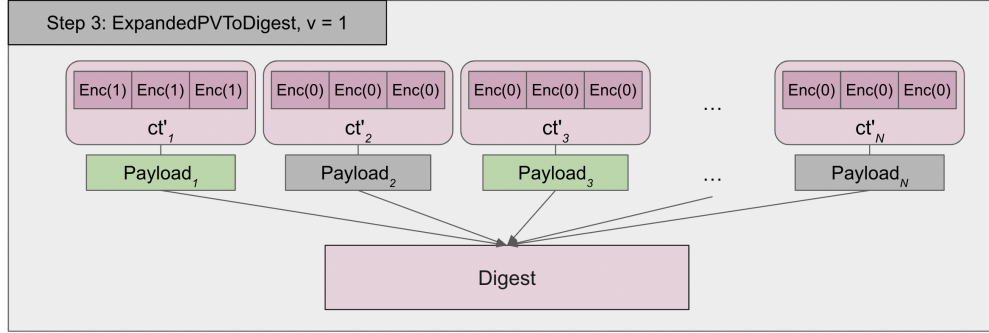
⁶Here $(\cdot)^D$ means a vector of D of \cdot elements.

⁷In [27], the authors provide an optimization to this step at the cost of a deeper circuit. We omit the details here for simplicity.

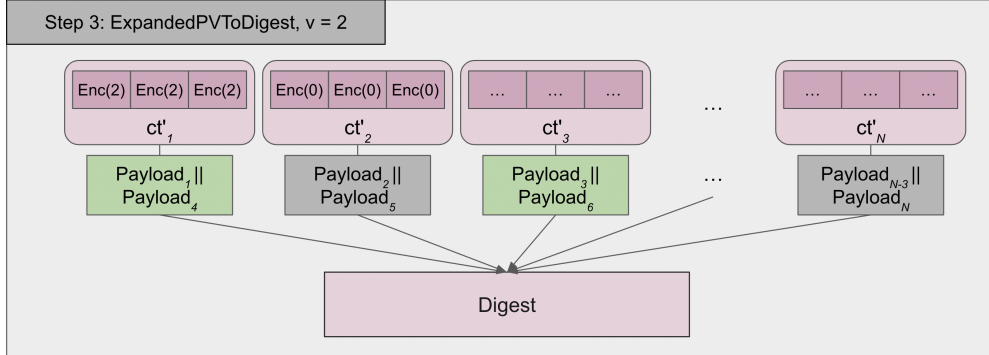
D multiplications and $D \log(D)$ rotations in total. Again, this realization focuses on $v = 1$ and later we show how larger v works.

4.2.3 Step 3: Use PV to Construct the Digest

Lastly, with all these N' ciphertexts above encrypting either non-zero (if pertinent) or zero (if impertinent), the detector forms a compact digest using BFV: if the i -th ciphertext does not encrypt zero, message i should be included in the digest. In OMRp2, $N' = N$, and each ciphertext either encrypts 1 or 0. However, our optimization requires those ciphertexts to encrypt $v > 1$ and $N' = N/v$. We visualize step 3 in Fig. 3. The interface is as follows:



(a) Step 3 with $v = 1$ (i.e., for the original OMR construction in Section 4.2.3)



(b) Step 3 with $v = 2$ (i.e., for our new construction in Section 5.2.3)

Figure 3: Visualization of step 3 ExpandedPVToDigest. When $v = 1$, each ciphertext has one corresponding payload to be included in the digest. When $v = 2$, since two messages are viewed as a single one, their payloads are concatenated.

- $M \leftarrow \text{ExpandedPVToDigest}(\text{pp}, \text{pk}_{\text{detect}}, (\text{ct}_1, \dots, \text{ct}_{N'}, \bar{k}), \text{BB})$: takes public parameter pp , a detector key $\text{pk}_{\text{detect}}$, a vector of BFV ciphertexts of length N' , a bulletin board BB of size $N \geq N'$, N divides N' , and an upper bound \bar{k} on the number of pertinent messages addressed to that recipient; outputs a digest M .

The correctness is as follows:

Definition 4.3 (Correctness of ExpandedPVToDigest). There exists a PPT algorithm `DecodeDigest` taking a digest M and a secret key sk and outputting payloads PL such that: for the same quantifiers as in Definition 4.1, for any $N' \leq N$, and any vector of ciphertexts $(\text{ct}_1, \dots, \text{ct}_{N'})$ encrypted under $\text{pk}_{\text{detect}} = \text{pk}_{\text{BFV}}$, let $M \leftarrow \text{ExpandedPVToDigest}(\text{pp}, \text{pk}_{\text{detect}}, (\text{ct}_1, \dots, \text{ct}_{N'}), \bar{k}, \text{BB})$ and $\text{PL} \leftarrow \text{DecodeDigest}(M, \text{sk})$; let $k = |S|$ (the number of pertinent messages in S), it holds that either $k > \bar{k}$ and $\text{PL} = \text{overflow}$, or:

$$\Pr[x_j \in \text{PL} \mid \text{BFV.Dec}(\text{sk}, \text{ct}_i) \neq 0^D] \geq (1 - \text{negl}(\lambda)) \quad \text{for all } j \in [N], i = j \bmod N' .$$

ExpandedPVToDigest Implementation (with $N' = N$). To realize `ExpandedPVToDigest`, [26] first encodes all the pertinent indices into the digest, and then the corresponding pertinent payloads.

We refer to the first part as “index encoding”: `OMRp2` first initialize $m > \bar{k}$ buckets, where \bar{k} is the upper bound of the number of pertinent messages. Then, it randomly assigns all the N messages into m buckets, and let Y_i , represent the set of messages (represented by indices) assigned to bucket $i \in [m]$.

For each bucket $i \in [m]$, compute $\text{Acc}_i \leftarrow \sum_{j \in Y_i} (\text{ct}_j \cdot j)$. If there is no pertinent message assigned to bucket i , Acc_i encrypts 0. If there is only one pertinent message j assigned to bucket i , Acc_i encrypts j , and the recipient can easily decrypt j to be the index of that pertinent message. However, if more than one pertinent message gets assigned in bucket i , there is a collision. To inform recipients of such collisions, `OMRp2` computes $\text{ctr}_i \leftarrow \sum_{j \in Y_i} \text{ct}_j$ for bucket i , which is the number of pertinent messages assigned to bucket i .⁸ The process is repeated $C \geq 1$ times to allow the recipient to obtain all the pertinent indices except with negligible probability, even with non-negligible probability of collision.

After obtaining all the pertinent indices, obtaining the pertinent payloads is easier. We refer to this second part as “payload encoding”. The detector first samples a uniform random matrix⁹ $A \in \mathbb{Z}_t^{K \times N}$ for some $K > \bar{k}$, and computes $\text{comb} \leftarrow (A \circ Z) \times (\text{ct}_1, \dots, \text{ct}_N)$, where $Z = \begin{pmatrix} x_1, x_2, \dots, x_N \\ \vdots \\ x_1, x_2, \dots, x_N \end{pmatrix} \in \mathcal{P}^{K \times N}$,¹⁰ for x_i being payload of message i (recall that \circ is the Hadamard product introduced in Section 3). With the pertinent indices and A (sent to the recipient as a random seed), the recipient uses Gaussian elimination to solve for all the pertinent payloads except with negligible probability.¹¹

The digest thus includes $B_j = ((\text{Acc}_{i,j})_{i \in [m]}, (\text{ctr}_{i,j}))$, $\forall j \in [C]$, together with comb , and the seed s used to generate A .

4.3 Decoding

The last step is for the recipient to run `Decode`. The recipient first checks all the counters $\text{ctr}_{i,j}$ ($i \in [m], j \in [C]$) and filters the ones that encrypt 1. The corresponding $\text{Acc}_{i,j}$ then contain all the pertinent indices. After obtaining these indices, the recipient uses the seed s to recover A and uses comb to solve for all the pertinent payloads.

⁸Note that if there are $t+1$ pertinent messages, there is an overflow since the calculation is done over \mathbb{Z}_t . However, [26] parametrizes the construction such that it overflows with negligible probability.

⁹Note that in [26, 27], a sparse matrix is used instead. However, that requires additional explanation and is not useful for us, so we ignore that part. See Remark 5.5 for more discussion.

¹⁰W.o.l.g., assume $\mathcal{P} := \{0,1\}^P$ can be embedded to \mathbb{Z}_t as if not, simply repeat this process for $P/\log(t)$ times.

¹¹As discussed in [26], $K = \bar{k} + \delta/\log(t)$ to achieve $O(2^{-\delta})$ failure probability.

5 Improved OMR Construction

This section focuses on our new construction of OMR, containing a new setup phase (a different way to generate the clue keys and the clue) and a more efficient algorithm for each of the three steps forming the detector retrieval.

5.1 Reducing the Clue Key Size using RLWE

We start with the setup algorithms and keys in Section 4.1. Recall that the setup for each recipient is essentially generating a PVW key pair, whose public part serves as the clue key. One major issue with the PVW scheme used in OMRp2 is that the public key size is $w\ell \log(q) = \omega(\ell n \log^2(q))$, and concretely, hundreds of kilobytes — which is awkward to distribute (e.g., it is too large for direct inclusion in a cryptocurrency wallet address). We introduce a tailored variant of RLWE encryption of [30] to resolve this issue.

The main contributor to the large PVW public key is the parameter w , which is the number of LWE samples it contains. $w = \omega(n \log(q))$ is required by the leftover-hash lemma to guarantee that $A\vec{e}, P\vec{e}$ are indistinguishable from uniformly random vectors over \mathbb{Z}_q .¹² Thus, we suggest an alternative strategy of encryption that avoids relying on the leftover hash lemma. Instead of making $A\vec{e}, P\vec{e}$ statistically indistinguishable from randomly drawn vectors, the scheme can rely on computational assumptions.

In other words, for $w = n$, although $A\vec{e}, P\vec{e}$ by themselves are not statistically close to random vectors, by adding some noise and get $A\vec{e} + \vec{x}', P\vec{e} + \vec{x}''$, where \vec{x}', \vec{x}'' are noise vectors, we again have the resulting public key indistinguishable from random vectors based on LWE assumption; and thus greatly reduces the public key size with $w = n$.

To achieve even better efficiency, instead of relying on LWE, we rely on RLWE. In more detail, the key generation algorithm samples $\alpha \xleftarrow{\$} \mathcal{R}_q$, where $\mathcal{R} := \mathbb{Z}(X)/(X^n + 1)$ for some security parameter n being a power-of-two. The secret key $s \in \mathcal{R}$ is sampled from some distribution \mathcal{D} , and the public key is, instead, $(\alpha, \beta = \alpha s + x)$ for some noise x sampled from noise distribution χ_σ . To encrypt, the sender simply samples $e \leftarrow \mathcal{D}$ and computes $a \leftarrow \alpha e + x', b \leftarrow \beta e + x''$ where $x', x'' \leftarrow \chi_\sigma$. Note that if we use a matrix $A \in \mathbb{Z}_q^{n \times n}$ to represent α , A is structured instead of being uniformly random from $\mathbb{Z}_q^{n \times n}$. Thus, using a ring element α is only possible given that we are not relying on the leftover hash lemma.

To make it more suitable for our use case, since we need to encrypt just $\ell \ll n$ bits, only the first ℓ coefficients of the ring element b are needed during decryption. In addition, to guarantee correctness with probability $1 - \epsilon_n$, the scheme simply needs to choose a range parameter r used for decryption to guarantee that the noise of the ciphertexts is $\leq r$ except with ϵ_n probability. With all noises sampled from χ_σ , and a distribution \mathcal{D} such that the Hamming weight of s, e drawn from \mathcal{D} are both bounded by h and $|s|_\infty, |e|_\infty = 1$, the aggregated noise of $(as + b)$ can be viewed as sampled from distribution $\chi_{\sqrt{2h+1}\cdot\sigma}$ (since there are in total $2h + 1$ independently sampled noise being summed up). We can thus set r according to $\chi_{\sqrt{2h+1}\cdot\sigma}$ to guarantee correctness.

Defining sRLWE encryption. Putting all these together, we get a tailored variant of the RLWE encryption formally stated as follows:

- $\text{pp} = (n, \ell, q, \sigma, r, \mathcal{D}) \leftarrow \text{sRLWE.GenParam}(1^\lambda, \ell, q, \sigma, \epsilon_n)$: Choose a secret key dimension n and a distribution \mathcal{D} where the distribution is sampling a random vector of form $\{-1, 0, 1\}^n$ with

¹²Recall that $A\vec{e}, P\vec{e}$ are computed during PVW.Enc for $A, P \in \mathbb{Z}_q^{n \times w} \times \mathbb{Z}_q^{n \times \ell}$ being the public key (Section 3.2).

a fixed Hamming weight h , such that $\text{RLWE}_{n,q,\mathcal{D},\sigma}$ holds. Set ciphertext modulus q , number of bits in plaintext $\ell \leq n$, and standard deviation σ for Gaussian distribution for ciphertext noise generation. Additionally, set minimum integer r such that $\text{erf}(\frac{r}{\sqrt{2} \cdot \sqrt{2h+1} \cdot \sigma}) \leq \epsilon_n / \ell$.

- $(\text{sk}, \text{pk}) \leftarrow \text{sRLWE.KeyGen}(\text{pp})$: Draw a secret key $s \leftarrow \mathcal{D}$. Sample $\alpha \xleftarrow{\$} \mathcal{R}_q$ and noise $x \leftarrow \chi_\sigma \in \mathcal{R}$, and compute $\text{pk} = (\alpha, \alpha s + x) \in \mathcal{R}_q \times \mathcal{R}_q$, $\text{sk} \leftarrow s$.
- $\text{ct} = (a, \vec{b}) \leftarrow \text{sRLWE.Enc}(\text{pp}, \text{pk}, \vec{m})$: To encrypt a vector $\vec{m} \in \mathbb{Z}_2^\ell$, define the ring element $t \leftarrow \sum_{i \in [\ell]} \frac{q}{2} \cdot \vec{m}[i] X^{i-1} \in \mathcal{R}_q$. Draw a ring element $e \leftarrow \mathcal{D} \in \mathcal{R}_q$ and noises $x', x'' \leftarrow \chi_\sigma$. Let $b \leftarrow t + x'' = \sum_{i \in [N]} b_i X^{i-1}$. The ciphertext is the pair $(a, \vec{b}) = (\alpha e + x', (b_i)_{i \in [\ell]}) \in \mathcal{R}_q \times \mathbb{Z}_q^\ell$.
- $\vec{m} \leftarrow \text{sRLWE.Dec}(\text{pp}, \text{sk}, \text{ct} = (a, \vec{b}))$: Let $a' \leftarrow a \cdot \text{sk} := \sum_{i \in [n]} a'_i X^{i-1}$, $\vec{d} = \vec{b} - (a'_i)_{i \in [\ell]}$, $\vec{m} = \lfloor \frac{\vec{d} + q/2}{r} \rfloor \in \mathbb{Z}_2^\ell$.

We now show that our construction has the same properties as the PVW encryption (Section 3.2), including the (tailored) correctness, CPA security, key privacy (i.e., ciphertexts under different public keys are computationally indistinguishable) and zero-plaintext wrong-key decryption (i.e., given the wrong key, a PVW ciphertext is decrypted into a non-zero plaintext with high probability), as follows.

Theorem 5.1. *Assuming the hardness of RLWE, sRLWE satisfies the following properties:*

- (Correctness) For any $\lambda > 0, q = \text{poly}(\lambda), \sigma > 0, 1 > \epsilon_n > 0$ and $\ell \leq n$ for n chosen in pp_{sRLWE} , let $\text{pp}_{\text{sRLWE}} \leftarrow \text{sRLWE.KeyGen}(1^\lambda, \ell, q, \sigma, \epsilon_n)$, $(\text{sk}, \text{pk}) \leftarrow \text{sRLWE.KeyGen}(\text{pp}_{\text{sRLWE}})$, for any message $\vec{m} \in \{0,1\}^\ell$, it holds that: $\Pr[\text{sRLWE.Dec}(\text{sk}, \text{sRLWE.Enc}(\text{pk}, \vec{m})) = \vec{m}] \geq 1 - \epsilon_n - \text{negl}(\lambda)$.
- (CPA security) For any PPT adversary $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$, for the same quantifiers above in correctness, let the adversary choose two messages $(\vec{m}_1, \vec{m}_2, \text{st}) \leftarrow \mathcal{A}_1(\text{pp}_{\text{sRLWE}}, \text{pk})$; let $b \xleftarrow{\$} \{1, 2\}$, $\text{ct} \leftarrow \text{sRLWE.Enc}(\text{pp}_{\text{sRLWE}}, \text{pk}, \vec{m}_b)$, it holds that: $|\Pr[\mathcal{A}_2(\text{st}, \text{ct}) = b] - \frac{1}{2}| \leq \text{negl}(\lambda)$.
- (Key privacy) For the same quantifiers above in CPA security, let $(\text{sk}', \text{pk}') \leftarrow \text{sRLWE.KeyGen}(\text{pp}_{\text{sRLWE}})$; let the adversary choose a message $(\vec{m}, \text{st}) \leftarrow \mathcal{A}_1(\text{pp}_{\text{sRLWE}}, \text{pk}, \text{pk}')$; let $\text{ct} \leftarrow \text{sRLWE.Enc}(\text{pp}_{\text{sRLWE}}, \text{pk}, \vec{m})$, $\text{ct}' \leftarrow \text{sRLWE.Enc}(\text{pp}_{\text{sRLWE}}, \text{pk}', \vec{m})$, it holds that: $|\Pr[\mathcal{A}_2(\text{st}, \text{ct}) = 1] - \Pr[\mathcal{A}_2(\text{st}, \text{ct}') = 1]| \leq \text{negl}(\lambda)$.
- (Zero-plaintext wrong-key decryption) For the same quantifiers above in correctness, let $(\text{sk}', \text{pk}') \leftarrow \text{sRLWE.KeyGen}(\text{pp}_{\text{sRLWE}})$; it holds that: $\Pr[\text{sRLWE.Dec}(\text{sk}', \text{sRLWE.Enc}(\text{pk}, 0^\ell)) = 0^\ell] \leq (\frac{r}{q})^\ell + \text{negl}(\lambda)$.

Proof. • (Correctness) The noise of a ciphertext comes from x, x', x'' all sampled from χ_σ . Since \mathcal{D} gives out a binary vector with hamming weight x , there are $2h$ Gaussian noise contributed by x, x' , sampled independently. Additionally, x'' is one additional Gaussian noise. In total, there are $2h + 1$ independently sampled noises from discrete Gaussian distribution $(0, \sigma)$. The resulting noise x_t is thus from discrete Gaussian distribution $(0, \sqrt{2h+1}\sigma)$. To satisfy $\Pr[r \geq x_t] \geq 1 - p$ for some probability $0 > p > 1$, it is required that $\text{erf}(\frac{r}{\sqrt{2} \cdot \sqrt{2h+1} \cdot \sigma}) \leq p$. As required by sRLWE.GenParam , $\text{erf}(\frac{r}{\sqrt{2} \cdot \sqrt{2h+1} \cdot \sigma}) \leq \epsilon_n / \ell$. Thus, by union bound, all $\ell \geq 1$ noises together are bounded by r with probability ϵ_n . Thus, the correctness property follows straightforwardly.

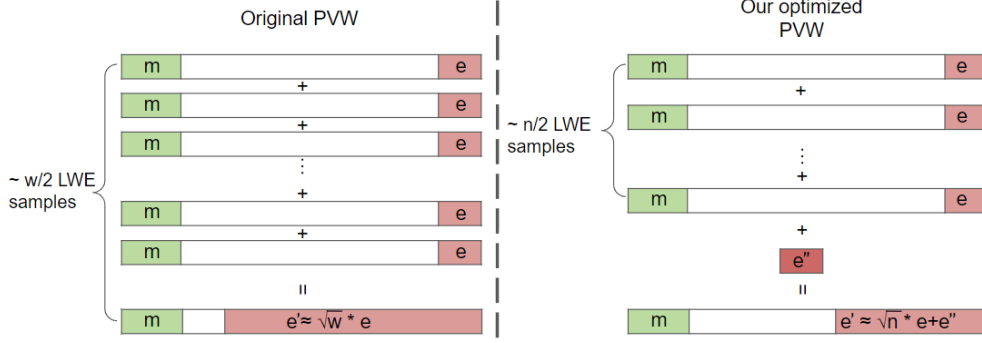


Figure 4: Visualization of the reduced noise of the ciphertexts from our RLWE encryption variant scheme compared to the original PVW scheme

- (CPA security) CPA security is the same as the original RLWE encryption in [30]. Then only change we make is that the ciphertext only contains $\vec{b} \in \mathbb{Z}_q^\ell$ instead of $b \in \mathcal{R}_q$. However, this does not affect the security guarantee (given that the RLWE assumption holds for our parameters).
- (Key privacy) We argue key privacy via a hybrid argument. We introduce the hybrid construction Π_1 : the only difference between Π_1 and sRLWE is that instead of computing $a \leftarrow \alpha e + x, b \leftarrow \beta e + x'' + t$ as in Enc (where $t \leftarrow \sum_{i \in [q]} \frac{q}{2} \cdot \vec{m}[i] X^{i-1} \in \mathcal{R}_q$ for \vec{m} being the input plaintext of ℓ bits), it first samples $a, b' \xleftarrow{\$} \mathcal{R}_q$ uniformly at random, and then compute $b \leftarrow b' + t$, and outputs (a, b') as the ciphertext. This is simple OTP and thus achieves key privacy trivially. Thus, if there exists an adversary that breaks sRLWE, since it cannot break Π_1 , it must break the Ring-LWE assumption.
- (Zero-plaintext wrong-key decryption) By RLWE, given a ciphertext $\text{ct} = (a, \vec{b})$ encrypted under sk' , \vec{b} is indistinguishable from a random vector sampled from \mathbb{Z}_q^ℓ with respect to sk sampled independently from sk' . Therefore, $\vec{b} - \text{ask}$ results in a random vector sampled from \mathbb{Z}_q^ℓ , which decrypts to 0^ℓ with probability $(\frac{2r+1}{q})^\ell + \text{negl}(\lambda)$. □

Efficiency. Since α can be represented as a random seed, and β has $n \log(q)$ bits, the public key size is now $n \cdot \log(q)$ bits, which is much smaller. The Enc runtime is also reduced from $\omega(n^2 \log(q))$ to $O(n \log(n))$ due to the smaller public key, as each ring element multiplication runs in $O(n \log(n))$ time using NTT.

Setup for our OMR construction. Our very first step is to replace the underlying PVW scheme in OMRp2 with our RLWE encryption variant sRLWE. This reduces the clue key and sender runtime.

Smaller range r . One additional property of sRLWE compared to PVW is that the value r is much smaller than in Section 3.2. sRLWE requires $\text{erf}(\frac{r}{\sqrt{2} \cdot \sqrt{2h+1} \cdot \sigma}) \leq \epsilon_n / \ell$, while the original PVW requires $\text{erf}(\frac{r}{\sqrt{2} \cdot \sqrt{w} \cdot \sigma}) \leq \epsilon_n / \ell$ as in [26]. Since $h = O(n)$ and $w = \omega(n \log(n))$, the new noise range is much smaller. We visualize the difference in Fig. 4.

5.2 A New Retrieval Circuit

Another major bottleneck of the prior OMR construction is the detector runtime. Therefore, we shift our focus to `Retrieve` and design a new and more efficient retrieval circuit.

5.2.1 A New Homomorphic Decryption Circuit for Step 1

Recall that in OMRp2 step 1 (Section 4.2.1), the detector homomorphically decrypts the PVW ciphertexts via a linear transformation followed by a degree- $(t - 1)$ polynomial (where t is the plaintext modulus of BFV, $t = q$ for q being the PVW ciphertext modulus). This means that the evaluation of the polynomial requires $O(t)$ homomorphic operations, which is costly for $t > 2D$, where D is the ring dimension of the underlying BFV.

Fortunately, since the new RLWE encryption variant has a much smaller range r (as discussed in Section 5.1 above and visualized in Fig. 4), a new more efficient circuit can be designed as follows. The fundamental goal is to compute the following function using a polynomial function:

$$f(x) = \begin{cases} 0 & \text{if } t - r \leq x \leq r \\ 1 & \text{o.w.} \end{cases} \quad (1)$$

. While it has t distinct points, implying a degree- $(t - 1)$ function, it can be evaluated more efficiently at the cost of having a higher degree, i.e.,

we represent $f(x) := (\prod_{i=-r}^r (x - i))^{t-1}$. This representation chiefly requires two steps of computation: $y = \prod_{i=-r}^r (x - i)$, checking whether $x \in [-r, r]$, if so, $y = 0$; otherwise, $y \neq 0$. By Fermat's Little Theorem, y^{t-1} returns 1 iff $y \neq 0$. Therefore, this representation is equivalent to Eq. (1). Furthermore, we optimize the evaluation of y to be $y = \prod_{i=0}^r (x^2 - i^2)$, which is equivalent as before, but has $r + 1$ multiplications instead of $2r$ multiplications, and thus can be evaluated more efficiently. Then again, $f(x) = y^{t-1}$.

Efficiency. In total, to evaluate $f(x)$, only $r + 1 + \log(t - 1)$ multiplications are needed.

This decrease in the number of multiplications comes at the cost of increasing the multiplicative degree from $\sim t$ to $\sim r \cdot t$. The overall effect of this tradeoff on running time depends on the parameters. To evaluate a function of degree k , each multiplication takes $O(D \text{polylog}(k))$ time. Therefore, our new representation takes $O((r + \log(t)) \cdot (\text{polylog}(r) + \text{polylog}(t)))$ time to evaluate, while the original degree- $(t - 1)$ polynomial needs $O(t \text{polylog}(t))$ time. As long as $r \ll t$, our method is more efficient. Concretely, with the parameters chosen in Section 7, our new representation can be evaluated $\sim 20x$ faster.

We formalize our new construction in Algorithm 1.

Theorem 5.2. *ClueToPackedPV in Algorithm 1 is correct (Definition 4.1) given the correctness of the underlying BFV scheme.*

Proof. Since the correctness of BFV is assumed, it remains to show the circuit is indeed the sRLWE decryption circuit. To compute two ring element multiplication, $c \leftarrow ab \in \mathcal{R}_q$ with ring dimension n being a power-of-two, let $c = \sum_{i=0}^{n-1} c_i X^i$, $a = \sum_{i=0}^{n-1} a_i X^i$, $b = \sum_{i=0}^{n-1} b_i X^i$, we have $c_i = \sum_{j=0}^i a_j \cdot b_{n-j} - \sum_{j=i+1}^n a_j \cdot b_{n-j}$. Therefore, it is straightforward that line 9 to line 10 computes the first ℓ -th coefficients of ask for ciphertext (a, \vec{b}) encrypted under sk . Then, together with line 12, they together directly compute the decryption circuit excluding the range check for \vec{m} at the end of sRLWE.Dec. Lastly, the range-check checks whether the result is in $[-r, r]$, if so, returns 0 and otherwise returns

Algorithm 1 Our new ClueToPackedPV

```
1: procedure PerfOMR1.ClueToPackedPV(pp, pkdetect = (pkBFV, ctsk, BB = (xi, ci)i∈[N])
2:                                     ▷ ctsk = Enc(skBFV, skpvw))
3:   Parse ci = (ai = ∑j ai,jXj-1,  $\vec{b}$  = (bi,l)l∈[ℓ])
4:   Let d ← ⌈N/D⌉
5:   for k ∈ [d] do                                     ▷ D is the ring dimension of the BFV scheme
6:     for l ∈ [ℓ] do
7:       for u ∈ [D] do
8:         Let  $\vec{a}_u = (\dots) \in \mathbb{Z}_q^{n \times 1}$ 
9:         Let A ← ( $\vec{a}_1 || \dots || \vec{a}_D$ ) ∈  $\mathbb{Z}_q^{N \times n}$ 
10:        Homomorphically compute ct1 ← skpvw × AT
11:        Let  $\vec{b}' \leftarrow b_{k \cdot d+1, l} || \dots || b_{k \cdot d+D, l}$ 
12:        Homomorphically compute ctk,l ←  $\vec{b}' - \text{ct}_1$ 
13:        BFV.Eval(pkBFV, ctvk,l, h ∘ g), where g(x) = ∏i=0r(x2 - i2) and h(x) = xt-1
14:        Homomorphically compute ctv ← ∏l∈[ℓ] ctk,l
15:   return (ctk)k∈[d]
```

1. As shown, $g(x)$ returns 0 if the input is within $[-r, r]$ and returns a non-zero value if not. Then $h(x)$ returns 0 if the input is 0 and returns 1 otherwise. Therefore, $f(x) = h(g(x))$ exactly computes the range check. Therefore, line 6 to lin 13 computes the sRLWE.Dec circuit. By the correctness and wrong-key decryption property proven for Theorem 5.1, we conclude the correctness of ClueToPackedPV. \square

5.2.2 An Efficient PV Unpacking Algorithm for Step 2

After step 1 (constructed above, defined in Definition 4.1), we obtain $d = \lceil N/D \rceil$ ciphertexts, each of which encrypts D bits, indicating whether the N messages are pertinent or not. For step 2 (Definition 4.2), we need to expand these into $N' = N/v$ ciphertexts, each of which encrypts a single integer in all of its D slots.

Recall that v is the bundle size (i.e., we bundle v messages into a single one for efficiency). We first focus on $v = 1$ as in OMRp2 (construction in prior work) for simplicity, and thus $N' = N$. OMRp2 achieves this by performing $O(N \log(D))$ homomorphic operations (or $O(\log(D))$ levels of multiplications but with $O(N)$ operations as in [27])¹³. In this section, we introduce an algorithm with $O(N)$ homomorphic operations and a single multiplication level.

Message Extraction. Let us start with a single ciphertext ct encrypting $(m_1, \dots, m_N) \in \mathbb{Z}_t^N$. Recall that in BFV, the message vector is first encoded into a polynomial before encryption (see Section 3.3). This encoding is to make the multiplications between messages easier over \mathbb{Z}_t , while our goal instead, is to unpack these messages into individual BFV ciphertexts. Thus, we first reverse this encoding and extract the message in each slot out to each coefficient of the encoded polynomial.

Specifically, recall that the encoding works as follows. The encoding scheme construct a polynomial $y(X) = \sum_{i \in [N]} y_i X^{i-1}$ with $m_i = y(\zeta_i)$, where ζ is the $2N$ -th primitive root of unity of t ,

¹³[27] gives a flexibility to trade off the number of operations vs. levels.

and $\zeta_i := \zeta^{3^i}$. Thus, a ciphertext ct encrypting (m_1, \dots, m_N) encrypts the polynomial $y(X)$.

Our first step is thus to revert this process: i.e., homomorphically change $y(X)$ to $m(X)$. This can be done by computing $\text{ct}' \leftarrow \text{ct} \cdot U^\top$, where

$$U := \begin{pmatrix} 1 & \zeta_0 & \zeta_0^2 & \cdots & \zeta_0^{N-1} \\ 1 & \zeta_1 & \zeta_1^2 & \cdots & \zeta_1^{N-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \zeta_{\frac{N}{2}-1} & \zeta_{\frac{N}{2}-1}^2 & \cdots & \zeta_{\frac{N}{2}-1}^{N-1} \\ 1 & \bar{\zeta}_0 & \bar{\zeta}_0^2 & \cdots & \bar{\zeta}_0^{N-1} \\ 1 & \bar{\zeta}_1 & \bar{\zeta}_1^2 & \cdots & \bar{\zeta}_1^{N-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \bar{\zeta}_{\frac{N}{2}-1} & \bar{\zeta}_{\frac{N}{2}-1}^2 & \cdots & \bar{\zeta}_{\frac{N}{2}-1}^{N-1} \end{pmatrix} \in \mathbb{Z}_t^{N \times N}$$

where $\bar{\zeta}_j := \zeta_j^{-1}$. The resulting ct' then encrypts the polynomial $m(X) = \sum_i m_i X^{i-1}$ (i.e., $\text{BFV.PartialDec}(\text{ct}') = m(X)$), as introduced in [28] denoted as `SlotToCoeff`.

Unpacking. After obtaining a ciphertext ct' encrypting a polynomial $m(X) = \sum_{i \in [D]} m_i X^{i-1}$, we want to obtain D ciphertexts $\text{ct}'_1, \dots, \text{ct}'_D$, such that each ciphertext ct'_i encrypts a constant polynomial $p_i(X) = m_i$ (recall that a constant polynomial encodes a vector $(m_i, \dots, m_i) \in \mathbb{Z}_t^D$, i.e., $p_i(\eta) = m_i$ for all $\eta \in \mathbb{Z}_t$). To do so, the detector performs the oblivious expansion procedure introduced in [3] and generalized by [2, 25]. This well-established procedure is recalled in Algorithm 2 `OExpand`.

Allowing $\mathbf{v} > 1$. Despite the great efficiency improvement with the unpacking technique above (both `SlotToCoeff` and `OExpand` take only $O(D)$ operations per ciphertext), fundamentally, this requires $O(N)$ homomorphic operations for N messages, which is still quite costly.

One natural idea is to bundle $\mathbf{v} \ll N$ messages as a single one (\mathbf{v} to be fixed later), reducing N messages to $N' = N/\mathbf{v}$ messages before performing this PV unpacking process. The number of operations thus reduces to $O(N')$. In more detail, with d input ciphertexts $(\text{ct}_1, \dots, \text{ct}_d)$ (assuming \mathbf{v} divides d for simplicity), the detector first divides them into \mathbf{v} chunks: $(\text{ct}_1, \dots, \text{ct}_{d/\mathbf{v}}), (\text{ct}_{d/\mathbf{v}+1}, \dots, \text{ct}_{2(d/\mathbf{v})}), \dots$. Then, it adds up all the chunks ciphertext-wise, i.e. computing $\tilde{\text{ct}}_i \leftarrow \sum_{j=0}^{\mathbf{v}-1} \text{ct}_{j \cdot (d/\mathbf{v}) + i}$ for $i \in [d/\mathbf{v}]$. This gives d/\mathbf{v} ciphertexts, each with D slots, where each slot encrypts the summation of \mathbf{v} slots of the input ciphertexts. After obtaining these d/\mathbf{v} ciphertexts, everything proceeds as the unpacking procedure described above (expanding each slot into a single ciphertext).

Putting everything together, we obtain our `PVUnpack` algorithm as in Algorithm 2.

Theorem 5.3. *PVUnpack in Algorithm 2 is correct (Definition 4.2) given the correctness of the underlying BFV scheme.*

Proof. For simplicity, we start with a single input ciphertext ct_1 . Recall that in Section 3.3, the ciphertext encrypts a polynomial $y(X) = \sum_{i \in [N]} y_i X^{i-1}$ and let $m_i \leftarrow y(\zeta_i)$. Let $\vec{y} = (y_i)_{i \in [N]}$, computing $\vec{m} \leftarrow \vec{y} \cdot U^\top$ gives $\vec{m}[i] = m_i$ for all $i \in [N]$ by [28]. Therefore, `tmp1` computed in Algorithm 2 gives a ciphertext encrypting $m(X) = \sum_{i \in [N]} m_i X^{i-1}$ assuming the correctness of BFV. Then, by the correctness of the oblivious expansion algorithm [3, Thm 1], let $(\text{ct}'_i)_{i \in [D]} \leftarrow \text{OExpand}(\text{tmp}_1)$, it holds that ct'_i encrypts a constant polynomial m_i . Therefore, we have $\text{Dec}(\text{sk}, \text{ct}'_i) = m_i^D$. This trivially generalizes to d ciphertexts (i.e., apply this step to each of the ct_i for $i \in [d]$). \square

Algorithm 2 Our new PVUnpack

```

1: procedure OExpand(ct) (Adapted from [3])
2:    $\triangleright$  All the keys needed to complete this procedure are assumed to be implicitly taken.
3:   res  $\leftarrow$  [ct]
4:   for  $i = 0$  to  $\log D$  do
5:     for  $j = 0$  to  $2^i - 1$  do
6:       tmp0  $\leftarrow$  res[ $j$ ]
7:       tmp1  $\leftarrow$  tmp0  $\cdot x^{-2^i}$ 
8:       tmp' $j$   $\leftarrow$  tmp0 + Substitute(tmp0,  $D/2^i + 1$ )
9:       tmp' $i+2^j$   $\leftarrow$  tmp1 + Substitute(tmp1,  $D/2^i + 1$ )
10:    res  $\leftarrow$  [tmp'0, ..., tmp' $2^{i+1}-1$ ]
11:   for  $i = 0$  to  $D$  do
12:     res[ $i$ ]  $\leftarrow$  BFV.Eval(res[ $i$ ],  $1/D, \times$ )
13:   return res
14: procedure SlotToCoeff(ct) (Adapted from [28])
15:

```

$$U := \begin{pmatrix} 1 & \zeta_0 & \zeta_0^2 & \dots & \zeta_0^{N-1} \\ 1 & \zeta_1 & \zeta_1^2 & \dots & \zeta_1^{N-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \zeta_{\frac{N}{2}-1} & \zeta_{\frac{N}{2}-1}^2 & \dots & \zeta_{\frac{N}{2}-1}^{N-1} \\ 1 & \bar{\zeta}_0 & \bar{\zeta}_0^2 & \dots & \bar{\zeta}_0^{N-1} \\ 1 & \bar{\zeta}_1 & \bar{\zeta}_1^2 & \dots & \bar{\zeta}_1^{N-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \bar{\zeta}_{\frac{N}{2}-1} & \bar{\zeta}_{\frac{N}{2}-1}^2 & \dots & \bar{\zeta}_{\frac{N}{2}-1}^{N-1} \end{pmatrix} \in \mathbb{Z}_t^{N \times N}$$

```

16:   Homomorphically compute res  $\leftarrow$  ct  $\cdot U^\top$ ,
17:   return res
18: procedure PerfOMR1.PVUnpack((ct1, ..., ct $d$ ), v)
19:   for  $i \in [d/v]$  do
20:     ct' $i$   $\leftarrow$   $\sum_{j \in [v]}$  ct $j \cdot v + i$ 
21:   for  $i = 0$  to  $d/v$  do
22:     tmp $i$   $\leftarrow$  SlotToCoeff(ct' $i$ )
23:     [ct' $i \cdot D + 1$ , ..., ct' $i \cdot D + D$ ]  $\leftarrow$  OExpand(tmp $i$ )
24:   return (ct' $i$ ) $i \in [N']$  (where  $N' = d \cdot D/v$ )

```

5.2.3 A General Encoding Procedure for Step 3

Lastly, we discuss digest encoding. As aforementioned, v messages are bundled as a single one for efficiency. Therefore, to deliver all the pertinent payloads, the payloads of these v messages bundled together should also be concatenated and viewed as a single large payload. In other words, we consider the concatenated payloads of the form $x'_i \leftarrow x_{i_1} || \dots || x_{i_v}$. The corresponding ciphertext ct'_i indicating whether x'_i is pertinent (if ct'_i encrypts a value > 1 , x'_i should be included in the digest and otherwise not)¹⁴. Hence, the detector performs digest encoding as if there were only N/v messages.

Overall, this speeds up the step 2 of the detector (nearly) v times. However, there are several issues to address:

Spurious payloads. The decoding algorithm will output all v messages in the bundle, some of

¹⁴Recall that ct'_i encrypts the summation of v slots in the input ciphertexts of PVUnpack corresponding to message x_{i_1}, \dots, x_{i_v} .

which may be impertinent. This is allowed by \mathbf{v} -soundness as defined in Definition 3.3, and often the application using OMR can filter the bundled impertinent messages as discussed in Remark 3.5.

Alternatively, we can do a final filtering (and thus attain the original OMR soundness definition [26, Thm 4.1]) by a small addition to the scheme: the sender will encrypt the payloads to the recipient, using an encryption scheme with the property that decryption using the wrong key is detectable. See Appendix A.1 for details.

Intra-bundle collisions. A more serious issue is that the index encoding process would fail if using the realization in Section 4.2.2. Recall that during the index encoding part of the original `ExpandedPVToDigest` step, each message is assigned to a bucket Y_j . Then, for each bucket Y_j , compute $\text{Acc}_j \leftarrow \sum_{i \in Y_j} i \cdot \text{ct}_i$ and $\text{ctr}_j \leftarrow \sum_{i \in Y_j} \text{ct}_i$. The decoding procedure views all the bucket j with ctr_j encrypting a value > 1 as having collisions and discards that bucket.

Now, when there is at most one pertinent message among each bundle of \mathbf{v} messages, ct'_i indeed encrypts 0 or 1 and the original decoding procedure (Section 4.2.3) works. However, when there are multiple pertinent messages (either true or false positives) in the same bundle, the counter ctr in the associate bucket (i.e., sum of ct'_i for the bundled message i in the bucket) encrypts a number greater than 1, causing those buckets to be treated as a collision of different (bundled) messages and discarded.

To solve this, we use a different index encoding scheme, supporting pertinency vector entries that are greater than 1. Let $N' \leftarrow \lceil N/\mathbf{v} \rceil$ be the total number of bundled messages. At a high level, we expand each single bit of $i \in [N']$ into $\log(\mathbf{v} + 1)$ bits and combine the expanded bits into a new index $i' \in [N' \cdot (\mathbf{v} + 1)]$. In this way, when multiplying i' by \mathbf{v} , no carrying occurs into the next expanded bit..

More formally, we process as follows. For any index $i \in [N']$, let $i[j]$ denotes the j -th bit of i . We define the following function to expand each bit into $\log(\mathbf{v} + 1)$ bits: $i' \leftarrow \text{BitExpand}(\mathbf{v}, i) := \sum_{j=1}^{\lceil \log(N'+1) \rceil} i[j](\mathbf{v}+1)^{j-1}$. In other words, this function represents a binary value using a $(\mathbf{v}+1)$ -ary value: $0, 1 \in \mathbb{Z}_2$ are encoded by $0, 1 \in \mathbb{Z}_{\mathbf{v}+1}$.

Then, the index encoding process uses the new index i' instead of i . In more detail, the detector randomly assigns each bundled message i into a bucket j for $j \in [m]$ (where m is the number of buckets and $m > \bar{k}$ for \bar{k} being the upper bound on the number of pertinent messages). Let Y_j represent the set of indices (of messages) assigned to bucket j . Then, as before, for bucket j , the detector computes $\text{Acc}_j \leftarrow \sum_{i \in Y_j} \text{ct}'_i \cdot i'$; and computes $\text{ctr}_j \leftarrow \sum_{i \in Y_j} \text{ct}'_i$ for bucket j .

The collision happens if and only if one of the following conditions happens: (1) $r_j > \mathbf{v}$ for $r_j \leftarrow \text{Dec}(\text{ctr}_j)$; (2) let $b_j \leftarrow \text{Dec}(\text{Acc}_j)$ and represent b_j using $(\mathbf{v}+1)$ -ary values, i.e. $b_j = \sum_i b_{j,i}(\mathbf{v}+1)^{i-1}$, there exists an i such that $b_{j,i} \neq 0 \wedge b_{j,i} \neq r_j$.

This process, again, is repeated for C times to guarantee that all the indices can be successfully recovered.

Toy example. We provide an example of this encoding method. Assume we have $N' = 32$ bundled messages, where bundle 15, 20, 25 contain pertinent messages, and $\mathbf{v} = 7$. After building the messages accordingly, we obtain $\text{ct}'_{15} = \text{Enc}(1)$, $\text{ct}'_{20} = \text{Enc}(5)$, $\text{ct}'_{25} = \text{Enc}(7)$, $\text{ct}'_{i \notin \{15, 20, 25\}} = \text{Enc}(0)$: We first get the binary representation of the indices: $15 = 001111$, $20 = 010100$, $25 = 011001$.

Then we encode all of them using `BitExpand`: $\text{BitExpand}(15, 7) = 000,000,001,001,001,001$, $\text{BitExpand}(20, 7) = 000,001,000,001,000,000$, and $\text{BitExpand}(25, 7) = 000,001,001,000,000,001$.

If there is no collision, i.e., if those three pertinent messages are assigned into different buckets, then these buckets' `Acc` encode $15 \cdot 1 = 000,000,001,001,001,001$, $20 \cdot 5 = 000,101,000,101,000,000$,

and $25 \cdot 7 = 000,111,000,111,000,000$. The three corresponding ctr's are $1 = 001, 5 = 101, 7 = 111$.

If there is a collision, say 15 and 20 collide, then **Acc** for that colliding bucket would be decrypted to $000,101,001,110,001,001$ while the counter is $1+5 = 6 = 110$; mismatch happens: $110 \neq 001 \neq 101$ (where $001,101$ appears in the **Acc** decryption above). The recipient can identify such a collision. On the other hand, if 15 and 25 collide, $\text{ctr} = 8 > v = 7$, which is clearly a collision.

Decoding. Recipient decoding is then straightforward: first checks whether the counter encrypts a number $> v$, and if so, there is a collision. Then, the recipient checks whether the **Acc** number matches the **ctr** number.

Again, after such index encoding, the payloads are encoded using the sparse matrix the same way as in Section 4.2.3, so we omit the details. We formalize all these in Algorithm 3.¹⁵

Theorem 5.4. *ExpandedPVTtoDigest in Algorithm 3 is correct (Definition 4.3 with DecodeDigest defined in Algorithm 3) given the correctness of the underlying BFV scheme.*

Proof. The correctness is satisfied if the following three conditions are satisfied, assuming no false positives (i.e., at most \bar{k} non-zero elements encrypted in the input ciphertexts):

1. Gaussian elimination succeeds (the linear combinations are linearly independent)
2. There are at most $t - 1$ pertinent messages assigned to each bucket (as the operations are homomorphically over \mathbb{Z}_t).
3. Index decoding and decoding is correct.

Satisfying condition (1) is guaranteed by line 9.

To satisfy condition (2), we need

$$\begin{aligned}
\Pr[X \geq t] &< \Pr[X \geq 2N/D] \quad (\text{since } N < Dt/2) \\
&= \Pr[X \geq 2(N/D)] \\
&= \Pr[X \geq 2(N/m)(m/D)] \\
&\leq \exp\left(-\frac{\delta^2}{2 + \delta} \frac{N}{d}\right) \quad (\text{by Chernoff bound, where } \delta = 2(m/D) - 1 = 2m' - 1, m := m/D) \\
&\leq \exp\left(-\frac{(2m' - 1)^2 t/2}{2m' + 1 d'}\right)
\end{aligned}$$

By using the union bound, the probability that none of the m buckets overflowing is $m \cdot \exp\left(-\frac{(2m' - 1)^2 t/2}{2m' + 1 d'}\right) \leq \text{negl}(\lambda)$, where $m' = O(\lambda)$. This is guaranteed by line 8. Note that it is okay that the **Acc** overflow, as it overflows only if **ctr** encrypts a number $> v + 1$ (by line 18 and the fact that **ctr** does not overflow) which is viewed as collision and thus discarded.

Lastly, for condition (3), we need to argue that our encoding/decoding scheme is indeed correct. To see this, we start with a single bit under plaintext (no homomorphic operation). Suppose for each bucket, the assigned message i has one corresponding bit b_i and a corresponding pertinency value $\rho_i \in [0, v]$. Then, the bucket computes $\text{Acc} \leftarrow \sum_i b_i \cdot \rho_i$ and $\text{ctr} \leftarrow \sum_i \rho_i$. Then, the recipient checks: if $\text{ctr} > v$, there is a collision; other wise, if $\text{ctr} \neq \text{Acc}$ and $\text{Acc} \neq 0$, there is a collision. If $\text{Acc} = \text{ctr} \leq v$, $b_i = 1$ for all i . Otherwise, if $\text{Acc} = 0$ and $\text{ctr} \leq v$, $b_i = 0$ for all i . Therefore, we can extend this process to $\log(N)$ bits. Since all the messages have their own indices, their indices differ by at least one bit. Therefore, if the checks go through, it straightforwardly implies that there is no collision. Otherwise, there is a collision.

¹⁵One may consider using the encoding scheme in [17]. In Appendix A.2, we discuss in detail why our solution is more efficient and specially tailored for our case.

Algorithm 3 Our new ExpandedPVTodigest

```

1: procedure PerfOMR1.ExpandedPVTodigest(pp, pkdetect, (ct1, . . . , ctN',  $\bar{k}$ ), BB = ((xi, ·))i∈[N])
2:    $\hat{k} \leftarrow \bar{k} + N \log(N)\epsilon_p$  ▷ Recall that  $\epsilon_p$  is the false positive rate included in pp.
3:   Choose  $C, m$  s.t:
4:   (1)  $C \cdot m$  is minimized
5:   (2) the index encoding fails with probability  $\text{negl}(\lambda)$  (i.e., decoding using eliminating the collisions fails with negligible probability)
6:   ▷ Failure probability is  $1 - \prod_{i=1}^{\hat{k}-1} (1 - (\frac{i}{m})^C)$  per [26, Sec 6.1.2]
7:   (3) each bucket is assigned at most  $t - 1$  messages except with negligible probability
8:   ▷ Overflow probability is  $m \cdot \exp(-\frac{(2m'-1)^2 t/2}{2m'+1 m'})$  where  $m' \leftarrow m/D$ 
9:   Choose  $K$  such that a random matrix in  $\mathbb{Z}_t^{K \times \bar{k}}$  is full rank with  $1 - \text{negl}(\lambda)$  probability
10:  ▷  $K = O(\bar{k} + \lambda/\log(t))$  as discussed in [26].
11:  for  $i \in [D], j \in [0, \frac{N'}{D} - 1]$  do
12:     $x'_{j \cdot D + i} \leftarrow x_{j \cdot v + 1} || \dots || x_{j \cdot v + v}$ 
13:  for  $i \in [C]$  do
14:    Initialize  $\text{Acc}_{i,u}, \text{ctr}_{i,u}$  for  $u \in [m]$ 
15:    for  $j \in [d/v]$  do
16:       $u \xleftarrow{\$} [m]$ 
17:       $\text{Acc}_{i,u} \leftarrow \text{Acc}_{i,j} + \text{ct}_j \cdot \text{BitExpand}(v, \text{binary}(j))$ 
18:      ▷ Note that for each accumulator  $\text{Acc}$ , the calculation needs to be split into multiple  $\mathbb{Z}_t$  elements. Each  $\mathbb{Z}_t$  element contain  $t' = \lfloor \log(t)/\log(v+1) \rfloor$  bits, and thus totally need  $\lceil \log(N/(v+1))/t' \rceil$   $\mathbb{Z}_t$  elements.
19:       $\text{ctr}_{i,u} \leftarrow \text{Acc}_{i,j} + \text{ct}_j$ 
20:       $A \xleftarrow{\$} \mathbb{Z}_t^{K \times N'}$ 
21:      Homomorphically computes  $\text{comb} \leftarrow (A \circ Z) \times (\text{ct}'_i)_{i \in [N']}$  where  $Z = \begin{pmatrix} x'_1, x'_2, \dots, x'_{N'} \\ \vdots \\ x'_1, x'_2, \dots, x'_{N'} \end{pmatrix} \in \mathcal{P}^{K \times N'}$ 
22:  return  $M = (s, (\text{Acc}_{i,u}, \text{ctr}_{i,u})_{i \in [C], u \in [m]}, \text{comb})$ 
23: procedure PerfOMR1.DecodeDigest( $M, \text{sk}$ )
24:    $\hat{k} \leftarrow \bar{k} + N \log(N)\epsilon_p$ 
25:   Parse  $M = (s, (\text{Acc}_{i,j}, \text{ctr}_{i,j})_{i \in [C], j \in [m]}, \text{comb})$ 
26:   Initialize an empty set  $P = \{\}$  to record all pertinent indices
27:   for  $i = 1$  to  $C$  do
28:     for  $j = 1$  to  $m$  do
29:        $a, c \leftarrow \text{BFV.Dec}(\text{sk}, (\text{Acc}_{i,j}, \text{ctr}_{i,j}))$ 
30:       If  $c > v$ , skip this iteration
31:       Parse  $a$  into  $(v+1)$ -ary numbers, i.e.,  $a = \sum_i a_i (v+1)^i$ 
32:       If any  $a_i \neq c$ , skip this iteration
33:       Let  $a' = \sum_i a'_i$  where  $a'_i = 1$  if  $a_i \neq 0$ , and  $a'_i = 0$  if  $a_i = 0$ .
34:       Add  $a'$  to  $P$ 
35:   If  $|P| > \hat{k}$ , PL = overflow and skip the next step
36:   Use  $P, \text{comb}, s$  and Gaussian elimination to solve for all payloads PL
37:   If failed, PL = overflow
38:  return PL

```

Lastly, Excessive false positives could break completeness. However, the probability of having more than $N \log(N) \epsilon_p$ false positives is negligible, by the correctness of `ClueToPackedPV`. We set $\hat{k} \leftarrow \bar{k} + N \log(N) \epsilon_p$. All the bounds above are set with respect to \hat{k} and the argument holds. \square

Remark 5.5. As mentioned in Section 4.2.3, we do not use sparse matrix A but instead a uniform random matrix (see line 9 in Algorithm 3) as the weights to compute the random linear combinations of the payloads. In [26], a sparse matrix is suggested to boost efficiency (as if a cell is 0, the corresponding multiplication can be skipped). However, in practice, this may not be the case. This is mainly because by the SIMD nature of BFV ciphertext, a single ciphertext can store $D \cdot \log(t)$ bits of information, and thus can store $W = D \cdot \log(t)/P$ (for $\mathcal{P} = \{0,1\}^P$) linear combinations. In this case, as long as one of the W corresponding weights is non-zero for a particular payload, the whole ciphertext needs to be multiplied. Therefore, practically, a sparse matrix does not reduce the number of multiplications, unless \mathcal{P} is large enough (e.g., such that $W = 1$). One can easily change this step to use a sparse matrix using the Sparse Random Linear Coding (SRLC) discussed in [26, Section 6]. To avoid extra complicity, we omit the details about SRLC.

5.3 Putting Everything Together

Putting everything above together yields a more efficient OMR construction. The pseudocode is presented in Algorithm 4.

Theorem 5.6. *The scheme `PerfOMR1` in Algorithm 4 is an OMR scheme (with ν -soundness) for $N < D \cdot t/2$, assuming the hardness of RLWE, the correctness of BFV leveled HE. Moreover `PerfOMR1` is also ν -compact.*

Proof. • (Correctness) Correctness is straightforward given the proven correctness of `ClueToPackedPV`,

`PVUnpack`, `ExpandedPVToDigest`. In more detail, given the correctness of `ClueToPackedPV`, for pertinent messages, we obtain an encryption of 1 with probability $1 - \epsilon_n$. By the correctness of `PVUnpack`, the encryption of 1 is expanded to a BFV ciphertext encryption a non-zero element in all of its slots. Lastly, by the correctness of `ExpandedPVToDigest`, if the input BFV ciphertext encrypts a non-zero element, the recipient can decode the corresponding payload.

- (ν -soundness) It is easy to see that essentially the decoding algorithm solves for $\hat{k} = \tilde{O}(\bar{k} + N \epsilon_p) \leq \kappa$ variables with κ variables. Each variable has size $\nu \cdot |\mathcal{P}|$. Therefore, the total size of the output of `OMRp2.Decode` is trivially bounded by $\tilde{O}(\nu \cdot |\mathcal{P}| \cdot (\bar{k} + N \epsilon_p))$.
- (Privacy) Privacy is directly implied by the key-privacy property of `sRLWE` proven given the hardness of RLWE.
- (ν -compactness) Compactness is also straightforward. By the correctness of `ClueToPackedPV`, a message is detected as pertinent in a false positive way with probability ϵ_p . Therefore, the total number of encryption of 1's from `ClueToPackedPV` is $\tilde{O}(\bar{k} + \epsilon_p N)$.

By the correctness of `PVUnpack` and `ExpandedPVToDigest`, the size of $(\text{Acc}_{i,u}, \text{ctr}_{i,u})_{i \in [C], u \in [m]}$ is therefore $\tilde{O}(\log(\nu) \cdot (\bar{k} + \epsilon_p N))$ as each value in the `Acc` is bounded by νN and the counter is bounded by ν . Then, we have $C \cdot m$ are bounded by $\tilde{O}(\hat{k})$ as proven in [26, Thm 6.2].

Then, the size of `comb` is $\tilde{O}(\nu \cdot |\mathcal{P}| \cdot \kappa)$. By [26, Lemma 6.5], we have $\kappa = \tilde{O}(\bar{k} + \epsilon_p N)$, $|\text{comb}|$ is $\tilde{O}(\nu \cdot |\mathcal{P}| \cdot (\bar{k} + \epsilon_p N))$ (as the concatenated payloads have size $\nu \cdot |\mathcal{P}|$).

\square

¹⁶Such a choice exists since, given a fixed q , we can choose σ such that $r = \text{poly}(\sigma)$ is small enough (e.g., $r = (q-1)/2$ and sufficiently large ℓ to satisfy the condition), then choose n sufficiently large to maintain the security level. Our

Algorithm 4 PerfOMR1: Practical Oblivious Message Retrieval

Let $f_s(x)$ be a PRF. Let BFV and sRLWE be as defined above.

- 1: **procedure** PerfOMR1.GenParam($1^\lambda, \epsilon_p, \epsilon_n$)
 - 2: Choose $\text{pp}_{\text{BFV}} = (D, t, \dots)$ such that homomorphically evaluate Retrieve with all but negligible probability
 - 3: Choose $\ell, q = t, \sigma$ such that with the r generated below, it satisfies that $(\frac{2r+1}{q})^\ell \leq \epsilon_p$ ¹⁶
 - 4: $\text{pp}_{\text{PVW}} = (n, \ell, q, \sigma, r, \mathcal{D}) \leftarrow \text{sRLWE.GenParam}(1^\lambda, \ell, q, \sigma, \epsilon_n)$
 - 5: **return** $\text{pp} = (\text{pp}_{\text{BFV}}, \text{pp}_{\text{PVW}}, \epsilon_p, \epsilon_n)$
 - 6: **procedure** PerfOMR1.KeyGen(pp)
 - 7: $(\text{sk}_{\text{sRLWE}}, \text{pk}_{\text{sRLWE}}) \leftarrow \text{sRLWE.KeyGen}(\text{pp}_{\text{BFV}})$
 - 8: $(\text{sk}_{\text{BFV}}, \text{pk}_{\text{BFV}}) \leftarrow \text{BFV.KeyGen}(\text{pp}_{\text{BFV}})$
 - 9: $\text{ct}_{\text{sk}_{\text{sRLWE}}} \leftarrow \text{BFV.Enc}(\text{pk}_{\text{BFV}}, \text{sk}_{\text{sRLWE}})$
 - 10: **return** $(\text{sk} = (\text{sk}_{\text{BFV}}), \text{pk} = (\text{pk}_{\text{clue}} = \text{pk}_{\text{sRLWE}}, \text{pk}_{\text{detect}} = (\text{pk}_{\text{BFV}}, \text{ct}_{\text{sk}_{\text{sRLWE}}}))$)
 - 11: **procedure** PerfOMR1.GenClue($\text{pp}, \text{pk}_{\text{clue}}, x$)
 - 12: $\vec{b} \leftarrow (0, \dots, 0) \in \mathbb{Z}_2^\ell$
 - 13: $c \leftarrow \text{sRLWE.Enc}(\text{pk}_{\text{clue}}, \vec{b})$
 - 14: **return** c $\triangleright c \in \mathcal{R}_q \times \mathbb{Z}_q^\ell$
 - 15: **procedure** PerfOMR1.Retrieve($\text{pp}, \text{BB}, \text{pk}_{\text{detect}}, \bar{k}$)
 - 16: Select v such that the runtime of ClueToPackedPV and PVUnpack are similar and also v divides N/D
 - 17: $(\text{ct}_i)_{i \in [N/D]} \leftarrow \text{ClueToPackedPV}(\text{pp}, \text{pk}_{\text{detect}}, \text{BB})$
 - 18: $(\text{ct}'_i)_{i \in [N']} \leftarrow \text{PVUnpack}(\text{pp}, \text{pk}_{\text{detect}}, (\text{ct}_i)_{i \in [N/D]}, v)$
 - 19: $M \leftarrow \text{ExpandedPVToDigest}(\text{pp} = (c, C, m), \text{pk}_{\text{detect}}, (\text{ct}'_1, \dots, \text{ct}'_{N'}, \bar{k}), \text{BB})$
 - 20: **return** M
 - 21: **procedure** PerfOMR1.Decode(M, sk)
 - 22: **return** PerfOMR1.DecodeDigest(M, sk)
-

5.4 Optimizations

Next, we introduce several implementation-level optimizations that further improve efficiency. Note that these optimizations are all implementation-level and compute the same algorithm proposed above, so all the properties trivially hold.

One-time rotation of sk . Originally in [26], OMRp2 uses PVW, and thus there are ℓ secret key vectors independently sampled (forming an sk matrix). During homomorphic decryption (i.e., step ClueToPackedPV), when computing the vector-matrix multiplication $a \cdot \text{sk}$, the ℓ secret key vectors need to be rotated for a total of $n \cdot \ell$ times (note that only a single encryption of sk and a single rotation key is sent to minimize the size of the detection key). However, with our use of sRLWE, we compute $a \cdot \text{sk}$ over \mathcal{R}_t , where sk is a single \mathbb{Z}_t^n vector and a is a (structured) matrix. This means that we only need to rotate sk n times. Furthermore, we rotate the sk once at the beginning of the retrieval and store them all in the memory. This implements $a \cdot \text{sk}$ for all N messages using just n rotations, instead of nN/D rotations for a naive implementation.

Improved linear transformation. Recall that in step PVUnpack, we need to compute a linear implementation uses a choice that minimizes the cost of homomorphically evaluating the induced decryption circuit.

transformation as in Algorithm 2 line 16, which takes N rotations and N homomorphic multiplications. We reduce this to just $2\sqrt{N}$ rotations using the Baby-step-giant-step algorithm of [20].

Two-level oblivious expansion. Recall that in `OExpand` in Algorithm 2, we expand a single ciphertext with D slots to D ciphertexts. However, a naive call of this function directly produces D ciphertexts. Concretely this is very costly: each ciphertext has a size $\sim 100\text{KB}$, so with $D = 32768$ (cf. Section 7), this takes 3.2GB of memory (and it gets worse for multi-thread). To reduce the memory cost, we optimize the `OExpand` to have two levels. We first expand the input ciphertext to l_1 ciphertexts (to-be-fixed later), and then, for each expanded ciphertext, further expand it to $l_2 = D/l_1$ ciphertexts. In this way, we keep only $l_1 + l_2$ ciphertexts in memory and expand on the fly for each batch of l_2 ciphertext. By setting $l_1 \approx \sqrt{D}$, this greatly reduces the memory usage.

6 OMR with Further Reduced Key and Clue Sizes

While `PerfOMR1` in Section 5 greatly improves the detector runtime, which is indeed a major practical concern of deploying OMR in real-world applications, there are some other practical considerations.

Recall that in the previous construction, the clue key size is reduced from $w \cdot \log(q)$ to $n \cdot \log(q)$ and the clue size remains to be $\sim(n \cdot \log(q))$. However, since now `sRLWE` relies on sparse secret keys, practically, n needs to be larger than in [26] (which relies on a secret key uniformly randomly sampled from \mathbb{Z}_q^n) for the same security level. Therefore, the concrete clue size increases (of course, the clue key is still largely reduced compared to the original OMR construction). Additionally, since the circuit for homomorphic decryption is deeper (the decryption function in `ClueToPackedPV` now has degree $r \cdot t$ compared to t in the original construction), the detection key size (mainly the BFV evaluation key, poly-logarithmically in the depth) is also larger. According to the original [26, Sec 10] parameters, the clue has size $\sim 1\text{KB}$, which is almost as large as the Zcash transaction size (1.3KB). Therefore, it might be hard to directly put the clue as part of the transaction. Additionally, the detection key size is $> 100\text{ MB}$. While it is only a one-time cost sent from the client to the server, further enlarging it can be a burden to light clients. Lastly, the clue key size is always better to be smaller.

Taking these three costs into consideration, we propose an alternative construction that makes these reduce the size of the clue, the detection key, and the clue key, at the cost of making the detection time larger than our protocol in Section 5, while still slightly smaller than the original construction in [26]. Our alternative construction provides a tradeoff between these metrics and the detector runtime, while still fully superseding the original construction in terms of efficiency.

6.1 Reducing the Clue Field

In `PerfOMR1` above (as in `OMRp2` [26]), we use $q = t$ so that decryption of the clues' PVW ciphertexts (defined over \mathbb{Z}_q) can be done directly via BFV homomorphic operations (defined over \mathbb{Z}_t). This modulus matching provides the mod q reductions for free. However, it forces q to match the large t needed by BFV, and we pay for this in the size of clues and clue keys (both $O(n \cdot \log(q))$).

Removing this constraint gives us the flexibility to reduce the clue size, and we indeed do so when `sRLWE` uses (sparse) short secrets, as follows.

We set $q \ll t$ (q to be even for simplicity) such that evaluating the PVW decryption in \mathbb{Z}_t is the same as working in \mathbb{Z} (i.e., no modular reductions). Then, when the PVW decryption performs

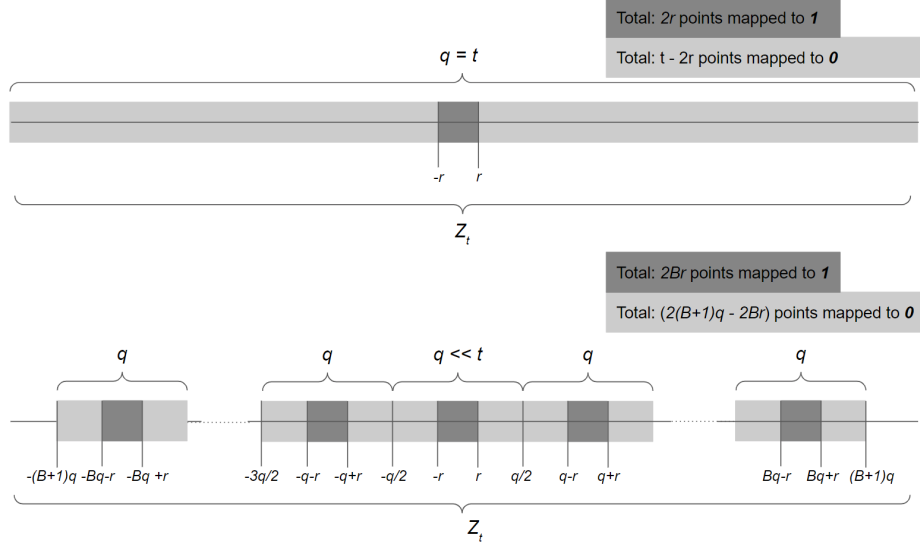


Figure 5: Visualization of the difference between setting $q = t$ and $q \ll t$ for sRLWE, where q is the sRLWE ciphertext modulus and t is the BFV plaintext modulus.

its \mathbb{Z}_q range test, we need to account for all the congruent values in \mathbb{Z} (see Fig. 5).

Formally, for clue $(a, \vec{b}) \in \mathcal{R}_q \times \mathbb{Z}_q^\ell$, by computing $\vec{b} - (a'_i)_{i \in [\ell]}$ where $a' = a \cdot s$ over \mathcal{R} (instead of \mathcal{R}_q), we obtain $(\mu_i + \kappa_i \cdot q)$ for some $\kappa_i \in \mathbb{Z}$ satisfying $\|\mu_i + \kappa_i \cdot q\| \leq t/2$ for all $i \in [\ell]$. If the message is pertinent, $\mu_i \in [-r, r]$ for the noise bound r . Otherwise, μ_i is indistinguishable from uniform in $[-q/2, q/2)$.

Additionally, since s has some fixed Hamming weight h , we can trivially bound $|\kappa| \leq h + 1$. However, to make our construction more efficient, we bound κ more tightly as follows: no matter whether a message is pertinent or not, $a \in \mathcal{R}_q$ is indistinguishable from a uniformly drawn element from \mathcal{R}_q : this is because $a \leftarrow \alpha e + x$ where e, x are both drawn independently from s by the sender, and thus by the hardness of RLWE, a is computationally indistinguishable from random. Therefore, we define $X := \sum_{i \in [h]} u_i$ where $u_i \stackrel{\$}{\leftarrow} [-q, q]$, and as shown above, $a'_0 \approx_c X$, and the expected value of X is simply 0. Therefore, we can set the bound B such that $\Pr[|X| \geq Bq] < e^{-\frac{B^2}{h}} \leq \text{negl}(\lambda)$ (by additive Chernoff bound). Since there are totally $N \cdot \ell$ a'_i 's and $N \cdot \ell = \text{poly}(\lambda)$, the probability that all of them are bounded by $[-Bq, Bq]$ is $1 - \text{negl}(\lambda)$. Note that additionally, there is a $b - a'_i$ operation after computing a' , and thus the bound of $b - a'_i$ is $(B + 1)q$.

A new homomorphic decryption circuit. The homomorphic decryption in \mathbb{Z}_t can thus be realized as follows: first compute $c_i \leftarrow \vec{b}[i] - a'_i \in \mathbb{Z}_t$ for all $i \in [\ell]$ as before; then homomorphically evaluate the function f' over c_i :

$$f'(x) = \begin{cases} 0 & \text{if } x \in [\kappa \cdot q - r, \kappa \cdot q + r] \text{ for all } \kappa \in [-B, B] \\ 1 & \text{else if } x \in [-(B + 1)q, (B + 1)q] \end{cases}$$

This function has $2(B + 1)q$ distinct points, and therefore degree $2(B + 1)q - 1$ (requiring $O(Bq)$ operations). We can evaluate it as a polynomial using Lagrange interpolation.

Efficiency. The clue key and clue are both of size $O(n \cdot \log(q))$, and thus shrink now as we choose a smaller q . Furthermore, since q is smaller, we can maintain the noise distribution while reducing

Algorithm 5 PerfOMR2: Practical Oblivious Message Retrieval

PerfOMR2.KeyGen, PerfOMR2.GenClue, and PerfOMR2.Decode are exactly the same as PerfOMR1.KeyGen, PerfOMR1.GenClue, and PerfOMR1.Decode, so we omit the details.

- 1: **procedure** PerfOMR2.GenParam($1^\lambda, \epsilon_p, \epsilon_n$)
 - 2: Choose $\text{pp}_{\text{BFV}} = (D, t, \dots)$ and $\text{pp}_{\text{PVW}} = (n, \ell, q, \sigma, r, \mathcal{D})$ as follows:
 - 3: (1) homomorphically evaluate **Retrieve** with all but negligible probability
 - 4: (2) \mathcal{D} is a distribution sampling a random vector from $\{0,1\}^n$ with a fixed hamming weight h
 - 5: (3) Let $X := \sum_{i \in [h]} u_i$ where $u_i \stackrel{\$}{\leftarrow} [-q, q]$, set the bound B to be $\Pr[|X| \geq Bq] \leq \text{negl}(\lambda)$, and $(B+1)q \leq t/2$
 - 6: (4) $\ell \leq n$, $\text{erf}\left(\frac{r}{\sqrt{2} \cdot \sqrt{2h+1} \cdot \sigma}\right) \leq \epsilon_n/\ell$, $\left(\frac{2r+1}{q}\right)^\ell \leq \epsilon_p$
 - 7: (5) $\text{RLWE}_{n,q,\mathcal{D},\sigma}$ holds
 - 8: Choose $\ell, q = t, \sigma$ such that with the r generated below, $\left(\frac{2r+1}{q}\right)^\ell \leq \epsilon_p$
 - 9: $\text{pp}_{\text{PVW}} = (n, \ell, q, \sigma, r, \mathcal{D}) \leftarrow \text{sRLWE.GenParam}(1^\lambda, \ell, q, \sigma, \epsilon_n/2)$
 - 10: **return** $\text{pp} = (\text{pp}_{\text{BFV}}, \text{pp}_{\text{PVW}}, \epsilon_p, \epsilon_n, B)$
 - 11: **procedure** PerfOMR2.ClueToPackedPV($\text{pp}, \text{pk}_{\text{detect}}, \text{BB}$)
 - 12: Same as PerfOMR1.ClueToPackedPV except that for line 13, replace it with $\text{BFV.Eval}(\text{pk}_{\text{BFV}}, \text{ct}_{v,l}, f')$ where f' is defined as follows:
$$f'(x) = \begin{cases} 0 & \text{if } x \in [\kappa \cdot q - r, \kappa \cdot q + r] \text{ for all } \kappa \in [-B, B] \\ 1 & \text{else if } x \in [-(B+1)q, (B+1)q] \end{cases}$$
 - 13: **procedure** PerfOMR2.Retrieve($\text{pp}, \text{BB}, \text{pk}_{\text{detect}}, \bar{k}$)
 - 14: Same as PerfOMR1.Retrieve except that line 17 is replaced with calling PerfOMR2.ClueToPackedPV.
-

n as well, which further reduces the sizes.

On the other hand, the degree of the function is now $2(B+1)q - 1$ instead of $r \cdot t - 1$. Thus, with careful parameter selection, the depth is greatly reduced as well. Therefore, the detection key size is also reduced.

The downside is that, evaluating f' homomorphically requires $O(Bq)$ homomorphic operations, compared to only $O(r + \log(t))$ operations before. Moreover, since the noise distribution remains the same, $\frac{2r+1}{q}$ becomes larger. Therefore, the guarantee the same wrong-key decryption probability (the false positive rate in OMR), ℓ needs to be larger. For the concrete results of this tradeoff, see Section 7.

This construction is formalized in Algorithm 5.

Theorem 6.1. *The scheme PerfOMR2 in Algorithm 5 is an OMR scheme (with v -soundness) for $N < D \cdot t/2$, assuming the hardness of RLWE, the correctness of BFV leveled HE. Moreover, OMRp2 is also v -compact.*

Proof. The only difference between Algorithm 5 and Algorithm 4 is that now we use f' instead of f in ClueToPackedPV. Therefore, we only need to argue that for any clue (a, \vec{b}) it holds that $\Pr[(\vec{b}[i] - (\text{ask}_{\text{pvw}})[i]) \in [-B \cdot q - r, B \cdot q + r] \forall i \in [\ell]] \leq \text{negl}(\lambda)$. This is simply guaranteed by line 5. \square

Extension to Group OMR. In the Group OMR setting, PerfOMR2 can be effectively combined with the technique of [27], as discussed in Appendix B.

7 Evaluation

7.1 Methodology

We implemented the above PerfOMR1 and PerfOMR2 schemes in a C++ library (with the optimizations introduced in Section 5.4) [36]. Our code extends the OMR library [35] and uses the SEAL [32] and PALISADE [38] libraries. We then compare our constructions to an improved version of OMRp2 introduced in [27]. The main improvement comes from PVUnpack and parameter selections. We run our implementation and the improved OMRp2 using Google Compute Cloud e2-standard-4 with 16GB RAM for a fair comparison.

Parameters. We chose the number of messages to be $N = 2^{19}$, and let $\bar{k} = k = 50$ (i.e., the upper bound and the real total number of pertinent messages are both 50)¹⁷.

For OMRp2, we reuse the parameters in [27] and guarantee ~ 115 -bit of computational security. Therefore, for the sRLWE used in PerfOMR1 we choose $n = 1024, q = 65537, \sigma = 0.5, h = 50, \ell = 2$. For the BFV used in PerfOMR1, we choose $D = 32768, Q \approx 2^{905}, t = q = 65537$.¹⁸ For the sRLWE used in PerfOMR2, we choose $n = 512, q = 400, \sigma = 0.5, h = 40, \ell = 6$. For the BFV used in PerfOMR2, we choose $D = 32768, Q \approx 2^{808}, t = 65537$. All these parameter settings guarantee > 128 -bit of security by [1]. Furthermore, same as in [26], we choose the false negative rate $\epsilon_n = 2^{-30}$ and false positive rate $\epsilon_p = 2^{-20}$, range $r = 32$ for PerfOMR1 and $r = 29$ for PerfOMR2, and the κ_i bound in Section 6 as $B = 29$. We choose $\ell_1 = 1024, \ell_2 = 32$ for our two-level oblivious expansion (see Section 5.4). We used $m = 400$ buckets and $C = 16$ trials. For the the random matrix A (line 9 in Algorithm 3), we choose $K = 53$.

Lastly, we choose $v = 8$ for PerfOMR1 and $v = 2$ for PerfOMR2 (discussed later).

7.2 Results

Representative costs. Table 2 summarizes the main cost metrics of all our schemes and the baseline (i.e., OMRp2 in [26] integrated with optimizations in [27]).

Our PerfOMR1 is about 13.8 faster than OMRp2 in terms of the detector runtime. Furthermore, the clue key size is reduced by roughly 60x. On the other hand, the clue size is about 2.2x worse (since sRLWE uses sparse keys, n for PerfOMR1 is larger than the n used for OMRp2). The detection key is also about 10% worse (due to the larger depth, we needed a larger Q). Lastly, since we choose $v = 8$, the digest size is also increased, but only by about 2.5x, since originally the digest contains 2 ciphertexts, one for index encoding and one for payload encoding. With the current parameter setting, we need 4 ciphertexts for payload encoding (the 8 concatenated payloads together have a

¹⁷For simplicity, since $k = \bar{k}$, we set $\hat{k} = \bar{k}$ instead of $\hat{k} = \bar{k} + N \log(N) \epsilon_p$. One can also alternatively view it as we choose $\bar{k} = 45$ and $\hat{k} = 50$.

¹⁸Originally, we used the incorrect parameters by accident, as pointed out by [24]. We modified and reported the correct parameters according to [24] and [1]. This parameter set slightly reduces the computational security of our scheme than originally reported (120 to 115), but has little effect on the runtime (only about 8% slower), which is why we chose it. Originally, we have $h = 31$ and r below being 19, but they are now fixed to $h = 48$ and $r = 32$. Note that these changes only affect the security and runtime of PerfOMR1.

	Detector Runtime (ms/msg)	Clue Key Size (kB)	Clue Size (Bytes)	Detector Key Size (MB)	Digest Size (Bytes/msg)	Recipient Runtime(ms)
OMRp2 [26, 27]	1 thread: 109.5 2 thread: 54.7 4 thread: 51.7	132.81	956	139	1.08	20
PerfOMR1 <i>Section 5</i>	1 thread: 7.93 2 thread: 3.95	2.13	2181	171	2.71	37
PerfOMR2 <i>Section 6</i>	1 thread: 39.64 2 thread: 19.84	0.56	583	140	1.08	20

Table 2: Comparison of cost metrics. Costs are per recipient. The bulletin contains $N = 2^{19}$ messages, of which $\bar{k} = k = 50$ are pertinent to the recipient. ms/msg and Bytes/msg are all amortized over N messages. Each message has 612 bytes of payload (as in [26, 27]).

size of 4896 bytes and one BFV ciphertext can encrypt at most 65536 bytes), but still only one ciphertext for index encoding. See below for how the runtime and digest size scale with v .

For PerfOMR2, the runtime is slightly better ($\sim 2.7x$) than OMRp2 (mainly due to the improved PVUnpack step for $v = 1$ in Section 4.2.2). The clue key size is drastically decreased: roughly 235x smaller than the clue key size of OMRp2. The clue is about 1.6x smaller. The detection key and digest size both remain roughly the same. Therefore, PerfOMR2 is strictly better than OMRp2, while having some complicated tradeoffs compared to PerfOMR1.

Improvement of each individual step. Fig. 6 shows the runtime breakdown for the three main steps of our schemes, compared to OMRp2 (all using $v = 1$ for fair comparison).

For PerfOMR1, our ClueToPackedPV is much faster than OMRp2 (about 13.8x faster); the bottleneck is PVUnpack, which is why we choose $v = 8$ for PerfOMR1 to optimize the runtime in the rest of our benchmarks. Conversely, for PerfOMR2, the runtime of ClueToPackedPV is roughly the same as OMRp2, which is thus the bottleneck. Therefore, we choose $v = 1$ for PerfOMR2. Moreover, our PVUnpack is about 5x faster than the PVUnpack in [26] even for $v = 1$.

Lastly, in both of our constructions, ExpandedPVToDigest remains similar to [26], while our encoding scheme requires slightly more operations.

How costs scale with v . As shown in Fig. 7, the costs of our two constructions changes with v : the runtime decreases as v increases. However, since it only boosts the PVUnpack step, enlarging v only works well when PVUnpack is a bottleneck. For example, for our PerfOMR1, the runtime with $v = 8$ is about 2x faster than the runtime with $v = 1$. Conversely, the runtime for PerfOMR2 is only about 20% faster when changing v from 1 to 8. Additionally, the digest size also grows with v linearly for both schemes (except for $v = 2$ which comes for free due to the SIMD nature of BFV, which can also improve the OMRp2 construction in [26, 27] but not by much).

Larger N and \bar{k} . To show scalability, we also test $N = 2^{21}$ and $N = 2^{23}$ for $\bar{k} = k = 50$; and fixing $N = 2^{19}$, we also test $\bar{k} = k = 100$ and $\bar{k} = k = 150$ (with $v = 8$ and 1 respectively for PerfOMR1, PerfOMR2 respectively as in Table 2).¹⁹ As shown in Table 3, the total runtime scales linearly with N , and grows slightly with larger \bar{k} . Digest size is essentially independent of N and scales linearly with \bar{k} . The scaling behavior is essentially the same as the prior construction OMRp2

¹⁹Note that for larger N , we need slightly larger noise budget. We set $Q \approx 2^{917}$ for simplicity of benchmarking. However, based on our experiments and estimation, $Q \sim 2^{910}$ is already enough, as asymptotically the noise grows with $\log(N)/2$.

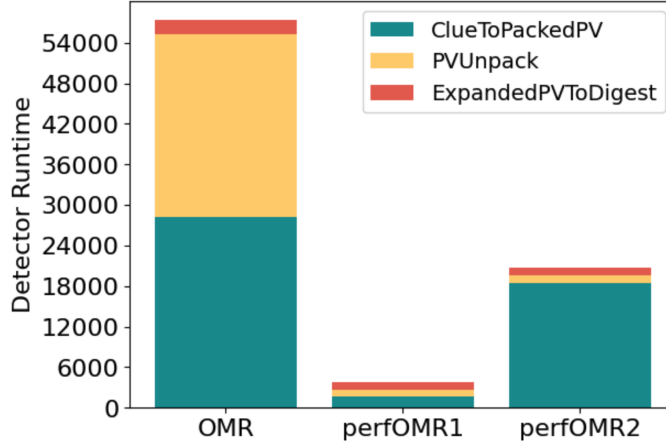


Figure 6: Comparison of runtime of each step for OMRp2 in [26], and our constructions PerfOMR1 and PerfOMR2 with $N = 2^{19}$ and $\bar{k} = k = 50$. We set $v = 1$ for both of our constructions for fair comparison.

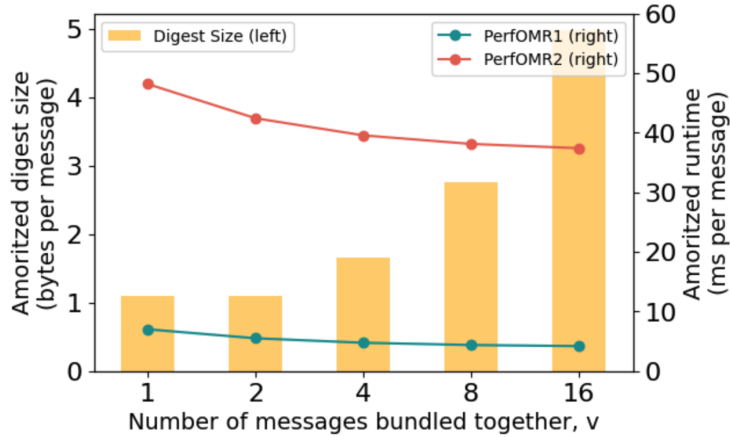


Figure 7: The runtime and digest size of our two schemes with respect to the value of v .

in [26, 27] and matches the asymptotic analysis.²⁰

Smaller ϵ_n, ϵ_p . To achieve even better ϵ_n, ϵ_p , we need to increase our parameters. For example, for $\epsilon_n = 2^{-80}$ and $\epsilon_p = 2^{-38}$ (which we believe is essentially enough for almost all real-world applications) for $\bar{k} = k = 45$ and $N = 2^{19}$, we need to make $r = 42$, $\ell = 4$, $C = 30$ for PerfOMR1 and other parameters remain unchanged.²¹ The runtime, by our estimation, is only about 2x slower and other metrics remain roughly the same based on our test. PerfOMR2 and OMRp2 in [26, 27] requires a similar parameter change and slow down.

²⁰As discussed in Remark 5.5, only when the payload size is large enough, does the runtime benefit from having sparse encoding and depends only on polylog of k . Of course, if k is too large, one can simply concatenate multiple payloads together to form a payload as large as 65536 bytes and take advantage of the sparse coding.

²¹Note that \bar{k} is reduced from 50 to 45 for simplicity. Otherwise, we need to increase K to 56 instead, which introduces another BFV ciphertext. This makes the runtime about 30% slower and the digest size 10% larger based on our estimation.

		$k = \bar{k} = 50$			
	N	Amortized runtime (ms/msg)	Total runtime (s)	Amortized digest size (Bytes/msg)	Total digest size (MB)
PerfOMR1	2^{19}	7.31	3931.65	2.71	1.35
	2^{21}		15868.37	0.68	
	2^{23}		64701.09	0.17	
PerfOMR2	2^{19}	39.64	20953.45	1.08	0.54
	2^{21}		82826.57	0.26	
	2^{23}		330985.56	0.06	

		$N = 2^{19}$			
	$k = \bar{k}$	Amortized runtime (ms/msg)	Total runtime (s)	Amortized digest size (Bytes/msg)	Total digest size (MB)
PerfOMR1	50	7.31	3931.65	2.71	1.35
	100	9.29	4874.03	4.87	2.43
	150	11.15	5847.55	7.04	3.52
PerfOMR2	50	39.96	20953.45	1.08	0.54
	100	41.58	21797.76	1.63	0.81
	150	42.85	22465.38	2.16	1.08

Table 3: Performance of our constructions when N and $k = \bar{k}$ varies.

7.3 Integration Considerations

Lastly, we discuss some system aspect of integrating the improved OMR in real-world applications, exemplified by the Zcash cryptocurrency [21] as analyzed in [26] (for detection cost, we consider Bitcoin-scale applications).

Clue key distribution. To integrate OMR, senders need to obtain the prospective recipient’s clue key to generate clues. As in [26], we consider Zcash’ Unified Addresses mechanism [22] to include a clue key as an extension of the recipient’s public address in a backward-compatible way, and extend the payment URIs similarly [34]. The clue key size of our PerfOMR1 is only 2.13KB, which can easily fit in a standard QR code that stores up to 3KB of data. Furthermore, our PerfOMR2 clue key size is only 0.59KB. Therefore, the clue key distribution is no longer as complicated as [26], which would have required some indirect mechanism (e.g., a URL) for fetching clue keys.

Clue embedding. For PerfOMR1, a clue of size 2181 bytes needs to be attached to every payload. This is larger than the 1.3kB of data on-chain per such payment. On the other hand, for PerfOMR2, the clue is only 583 bytes, much smaller than the one in OMRp2. This allows smaller on-chain data size. Including the clue in a transaction, we can simply use the `OP_RETURN` data that Zcash supports as discussed in [26].

Detection cost. The computational cost of detectors for OMRp2 is roughly \$1.95 per million payments scanned (4-thread), using commodity cloud computing.²² On the other hand, using our PerfOMR1, the cost is only \$0.12 per million payments scanned (2-thread). Using our PerfOMR2,

²²Using GCP `e2-standard-4`, 4 vCPUs, billed at \$0.136/hour.

the cost is \$0.88 per million payments scanned (2-thread). For Bitcoin-scale applications of roughly 300,000 payments per day,²³ our costs are \$1.12/month and \$7.88/month respectively, while prior work's cost is \$17.56/month.

Acknowledgements

We are grateful to the anonymous reviewers for their insightful comments and suggestions. We also thank Keewoo Lee and Yongdong Yeo for pointing out the parameter issues in Section 7.

This material is based upon work supported by DARPA under Contract No. HR001120C0085. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Defense Advanced Research Projects Agency (DARPA).

²³https://ycharts.com/indicators/bitcoin_transactions_per_day, retrieved 2023-10-17.

References

- [1] M. R. Albrecht, R. Player, and S. Scott. On the concrete hardness of learning with errors. *Journal of Mathematical Cryptology*, pages 169–203, 2015.
- [2] A. Ali, T. Lepoint, S. Patel, M. Raykova, P. Schoppmann, K. Seth, and K. Yeo. Communication–Computation trade-offs in PIR. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1811–1828. USENIX Association, Aug. 2021.
- [3] S. Angel, H. Chen, K. Laine, and S. T. V. Setty. PIR with compressed queries and amortized query processing. In *2018 IEEE S&P*. IEEE Computer Society Press, 2018.
- [4] G. Beck, J. Len, I. Miers, and M. Green. Fuzzy message detection. ACM CCS 2021.
- [5] E. Ben Sasson, A. Chiesa, C. Garman, M. Green, I. Miers, E. Tromer, and M. Virza. Zerocash: Decentralized anonymous payments from bitcoin. In *2014 IEEE S&P*, pages 459–474, 2014.
- [6] D. J. Bernstein. The chacha family of stream ciphers. <https://cr.yp.to/chacha.html>. Chacha.
- [7] J. Bethencourt, D. X. Song, and B. Waters. New techniques for private stream searching. *ACM Trans. Inf. Syst. Secur.*, 12:16:1–16:32, 2009.
- [8] A. Bittau, Ú. Erlingsson, P. Maniatis, I. Mironov, A. Raghunathan, D. Lie, M. Rudominer, U. Kode, J. Tinnes, and B. Seefeld. Prochlo: Strong privacy for analytics in the crowd. In *SOSP*, pages 441–459, 2017.
- [9] S. Bowe, A. Chiesa, M. Green, I. Miers, P. Mishra, and H. Wu. Zeze: Enabling decentralized private computation. In *2020 IEEE S&P (SP)*, pages 947–964, 2020.
- [10] Z. Brakerski. Fully homomorphic encryption without modulus switching from classical GapSVP. In *CRYPTO 2012*, LNCS. Springer, Aug. 19–23, 2012.
- [11] B. Chor, N. Gilboa, and M. Naor. Private information retrieval by keywords, 1998.
- [12] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan. Private information retrieval. In *36th FOCS*, pages 41–50. IEEE Computer Society Press, Oct. 23–25, 1995.
- [13] H. Corrigan-Gibbs, D. Boneh, and D. Mazières. Riposte: An anonymous messaging system handling millions of users. In *2015 IEEE S&P*, pages 321–338, 2015.
- [14] G. Danezis and C. Diaz. Space-efficient private search with applications to rateless codes. In *FC’07*, page 148–162. Springer, 2007.
- [15] J. Fan and F. Vercauteren. Somewhat practical fully homomorphic encryption. Cryptology ePrint Archive, Report 2012/144, 2012. <https://ia.cr/2012/144>.
- [16] M. Finiasz and K. Ramchandran. Private stream search at almost the same communication cost as a regular search. In L. R. Knudsen and H. Wu, editors, *Selected Areas in Cryptography*, pages 372–389, Berlin, Heidelberg, 2013. Springer.

- [17] N. Fleischhacker, K. G. Larsen, and M. Simkin. How to compress encrypted data. In C. Hazay and M. Stam, editors, *Advances in Cryptology – EUROCRYPT 2023*, pages 551–577, Cham, 2023. Springer Nature Switzerland.
- [18] R. Geelen and F. Vercauteren. Bootstrapping for bgv and bfv revisited. *J. Cryptol.*, 36(2), mar 2023.
- [19] C. Gentry. Fully homomorphic encryption using ideal lattices. In *ACM Symposium on Theory of Computing*, STOC '09, page 169–178. ACM, 2009.
- [20] S. Halevi and V. Shoup. Design and implementation of HElib: a homomorphic encryption library. Cryptology ePrint Archive, Report 2020/1481, 2020. <https://eprint.iacr.org/2020/1481>.
- [21] D. Hopwood, S. Bowe, T. Hornby, and N. Wilcox. Zcash Protocol Specification Version 2021.2.14. <https://github.com/zcash/zips/blob/master/protocol/protocol.pdf>.
- [22] D. Hopwood, N. Wilcox, T. Hornby, J. Grigg, S. Bowe, K. Nuttycombe, and Y. T. Lai. Zcash improvement proposal 316: Unified addresses and unified viewing keys. <https://zips.z.cash/zip-0316>, Apr. 2021.
- [23] S. Jakkamsetti, Z. Liu, and V. Madathil. Scalable private signaling. Cryptology ePrint Archive, Paper 2023/572, 2023. <https://eprint.iacr.org/2023/572>.
- [24] K. Lee and Y. Yeo. SophOMR: Improved oblivious message retrieval from SIMD-aware homomorphic compression. Cryptology ePrint Archive, Paper 2024/1814, 2024.
- [25] C. Lin, Z. Liu, and T. Malkin. XSPIR: Efficient symmetrically private information retrieval from ring-LWE. In *Computer Security – ESORICS 2022: 27th European Symposium on Research in Computer Security, Copenhagen, Denmark, September 26–30, 2022, Proceedings, Part I*, page 217–236, Berlin, Heidelberg, 2022. Springer-Verlag.
- [26] Z. Liu and E. Tromer. Oblivious message retrieval. In Y. Dodis and T. Shrimpton, editors, *Advances in Cryptology – CRYPTO 2022*, pages 753–783, Cham, 2022. Springer Nature Switzerland. Full version: Cryptology ePrint Archive 2021; internal citations follow the latter’s numbering.
- [27] Z. Liu, E. Tromer, and Y. Wang. Group oblivious message retrieval. *S&P 2024. Full version on eprint* <https://ia.cr/2023/534>, 2023.
- [28] Z. Liu and Y. Wang. Amortized functional bootstrapping in less than 7ms, with $\tilde{O}(1)$ polynomial multiplications. Asiacrypt 2023. <https://eprint.iacr.org/2023/910>.
- [29] J. Lund. Technology preview: Sealed sender for signal. <https://signal.org/blog/sealed-sender/>, Oct. 2018.
- [30] V. Lyubashevsky, C. Peikert, and O. Regev. On ideal lattices and learning with errors over rings. *J. ACM*, 2013.
- [31] V. Madathil, A. Scafuro, I. A. Seres, O. Shlomovits, and D. Varlakov. Private signaling. USENIX Security 2022, 2022.

- [32] Microsoft SEAL (release 3.6). <https://github.com/Microsoft/SEAL>, Nov. 2020. Microsoft Research, Redmond, WA.
- [33] S. Noether. Ring signature confidential transactions for monero. Cryptology ePrint Archive, Paper 2015/1098, 2015. <https://eprint.iacr.org/2015/1098>.
- [34] K. Nuttycombe and D. Hopwood. Zcash improvement proposal 321: Payment request URIs. <https://zips.z.cash/zip-0321>, Aug. 2020.
- [35] Oblivious message retrieval implementation. <https://github.com/ZeyuThomasLiu/ObliviousMessageRetrieval>, Dec. 2021.
- [36] Performr implementation. <https://github.com/ObliviousMessageRetrieval/ObliviousMessageRetrieval/tree/performr>, Mar. 2024.
- [37] R. Ostrovsky and W. E. Skeith. Private searching on streaming data. In *CRYPTO*, 2005.
- [38] PALISADE lattice cryptography library (release 11.2). <https://palisade-crypto.org/>, June 2021.
- [39] C. Peikert, V. Vaikuntanathan, and B. Waters. A framework for efficient and composable oblivious transfer. In *CRYPTO 2008*, pages 554–571. Springer, 2008.
- [40] R. Player. *Parameter selection in lattice-based cryptography*. PhD thesis, Royal Holloway, University of London, 2018.
- [41] S. Pu, S. A. Thyagarajan, N. Döttling, and L. Hanzlik. Post quantum fuzzy stealth signatures and applications. *CCS 23*, 2023.
- [42] R. L. Rivest, L. Adleman, and M. L. Dertouzos. On data banks and privacy homomorphisms. *Foundations of Secure Computation*, pages 169–179, 1978.
- [43] I. A. Seres, B. Pejó, and P. Burcsi. The effect of false positives: Why fuzzy message detection leads to fuzzy privacy guarantees? In I. Eyal and J. Garay, editors, *Financial Cryptography and Data Security*, pages 123–148, Cham, 2022. Springer International Publishing.
- [44] J. Szefer. Survey of microarchitectural side and covert channels, attacks, and defenses. *Journal of Hardware and Systems Security*, 3:219–234, Sept. 2019.
- [45] D. I. Wolinsky, H. Corrigan-Gibbs, B. Ford, and A. Johnson. Dissent in numbers: Making strong anonymity scale. In *OSDI 12*, pages 179–182. USENIX, Oct. 2012.

A Additional discussion

A.1 Boosting Soundness

As mentioned in Section 5.2.3, there is a simple way to boost v -soundness to full soundness.

Recall that the issue was that the bundled payloads may contain impertinent payloads, and the recipients need to distill the ones that are not pertinent to them.

To achieve this, we can use a key-private symmetric key encryption scheme (e.g., Chacha20 [6]) to encrypt the payload. In other words, for each message x , instead of publishing it directly onto the bulletin board, the sender will first get $x' \leftarrow \text{SKE.Enc}(\text{sk}, x)$ and publish x' . Therefore, the digest would be decrypted into a bundle of x' s on the recipient's side. To guarantee that only pertinent messages get successfully retrieved, we also need this secret key encryption algorithm to have the wrong-key detection property, such that $\perp \leftarrow \text{SKE.Dec}(\text{sk}', \text{SKE.Enc}(\text{sk}, x))$, for $\text{sk}' \neq \text{sk}$. Note that the secret key sk used for encryption is derived via a non-interactive-key-exchange scheme via the public key of the recipient. The sender also includes its ephemeral public key so that the recipient can later use it for key exchange. After obtaining all the payloads during decoding, the recipient simply performs a decryption using the secret key (derived using the corresponding ephemeral public key). The wrong-key detection property allows the recipient to tell which messages in the bundle are for them and which ones are not. This method slightly enlarges the payload (by the size of the public key for key-exchange). However, the overall performance is almost the same.

A.2 An Alternative Way to do Index Encoding

One may also use [17] to perform the index encoding for Section 5.2.3. This work focuses on how to compress sparse encrypted data, as required by our index encoding process. However, they only provide solutions to two different scenarios. The first scenario is when the pertinency vector is encrypting 0/1, which does not apply to us when $v > 1$. The second scenario is much more general: they do not assume a pertinency vector at all. However, in this case, they need to sample a random string of size $O(\lambda)$ for each message, and use it to detect collision, which means that the size of each bucket is $O(\log(N) + \lambda)$. This makes it less efficient than our solution, where only $O(\log(N) + \log(v))$ bits are needed (given $v \ll N = \text{poly}(\lambda)$). Therefore, we do not use their technique for efficiency concerns.

B Extension to group setting

[27] extends the OMR setting to a group setting. The sender selects a group of G recipients to which it wants to send. Then, it uses the G clue keys to generate a clue. The message is detected as pertinent with high probability to all the recipients included in the group and is detected as impertinent with high probability to all the other recipients.

A naive solution is simply to use an OMR construction to generate G clues for the G recipients and the detector simply scans all the G clues and encodes the digests accordingly. However, the issue is that all three steps `ClueToPackedPV`, `PVUnpack`, `ExpandedPVToDigest` need to be repeated for G times. This is very costly.

The main technique in [27] is first to let each recipient hold an identity id . During clue generation, the senders construct a (linear) function f which takes an id such that $f(\text{id})$ returns the PVW ciphertext of the intended recipient and includes this function as the clue. The function is

designed such that the homomorphic evaluation is very efficient. After homomorphically evaluating this function, the detector simply needs to perform all three steps of retrieval for one time instead of G times, thus greatly reducing the runtime.

Extending PerfOMR2 to group setting. PerfOMR2 is very compatible with this technique, and they can be beneficially combined to reduce the detector runtime. The main reason is that our retrieval runtime is comparable to the OMR construction in [26, 27]. Therefore, as shown in [27], replacing the expensive OMR retrieval process over every individual clue with some cheap homomorphic linear transformations (i.e., evaluating f) can greatly reduce the overall runtime in the group setting. Moreover, the small clue size of PerfOMR2 means that the function f (when represented using coefficients) is also smaller, resulting in smaller clues in the group setting too.

Extending PerfOMR1 to group setting. On the other hand, PerfOMR1 is not very compatible with this technique, since it has a very fast ClueToPackedPV step. We can alternatively first perform ClueToPackedPV for G clues as in the naive solution, Then, we sum up the results from ClueToPackedPV over the G results. Then, PVUnpack and ExpandedPVToDigest can be performed only once. Of course, similarly, we need to use our new ExpandedPVToDigest encoding scheme as the summation can result in a value > 1 thus causing a similar issue as we mentioned in Section 5.2.3. Therefore, for G that is small (e.g., $G < 20$), this simpler extension has a better runtime than using the technique in [27]. Of course, with larger groups, the technique in [27] is still more favorable.

Since it is not the main focus of our work, we leave exploring how to extend our construction more efficiently to the group setting in detail to future works.