

# Fast White-Box Adversarial Streaming Without a Random Oracle

**Abstract.** We study problems in the white-box adversarially robust streaming model. Streaming algorithms work with a large stream of data and aim to respond to queries regarding desired statistics of the data using limited memory and fast per item processing time, i.e., update time. Recently, the question of adversarially robust streaming, where the stream is allowed to depend on the randomness of the streaming algorithm, has gained a lot of attention. In this work, we consider a strong white-box adversarial model (Ajtai et al. PODS 2022), in which the adversary has access to all past random coins and the parameters used by the streaming algorithm. We also consider and define a related white-box adversarially robust distributed model. We note that both notions are also related to property preserving hashing in the cryptography literature. We consider the sparse recovery problem, which has been used for tasks such as distinct element estimation, low rank approximation of matrices and tensors, and finding a maximum matching. The main drawback of all previous work is that it requires a *random oracle*, which is especially problematic in the streaming model since the amount of randomness is counted in the space complexity of a streaming algorithm. Another issue with previous constructions is that many of them suffer from either large update time or rely on non-standard cryptographic assumptions. We construct a near-optimal solution for the sparse recovery problem in white-box adversarial streams, based on the more standard subexponentially secure Learning with Errors assumption. Importantly, our solution is the first without a random oracle and with polylogarithmic per item processing time. We also give related results in the distributed model. Our constructions are based on homomorphic encryption schemes satisfying very mild structural properties that are currently satisfied by most known schemes.

## 1 Introduction

The streaming model of computation has emerged as an increasingly popular paradigm for the analysis of massive datasets, where the overwhelming size of the data places restrictions on the memory, computation time, and other resources available to the algorithm. In the streaming model, the input data is implicitly defined through a stream of data elements that sequentially arrive one-by-one. One can also delete a previous occurrence of an item. An algorithm makes one pass over the stream and uses limited space to approximate a desired function of the input. In the formalization of Alon, Matias, and Szegedy [5], the data is represented as an underlying  $n$ -dimensional vector that is initialized to  $0^n$ , and then undergoes a sequence of additive updates to its coordinates. The algorithm

aims to optimize the space and time to process these updates, while being able to respond to queries about this underlying vector.

This model captures key resource requirements of algorithms for many practical settings such as Internet routers and traffic logs [51,74,66,70,21,22], databases [62,25,29,52,64,16], sensor networks [40,31], financial transaction data [3], and scientific data streams [20,67]. See [57] for an overview of streaming algorithms and their applications.

A large body of work on streaming algorithms has been designed for *oblivious* streams, for which the sequence of updates may be chosen adversarially, but it is chosen independently of the randomness of the streaming algorithm. Recently, there is a growing body of work on robust streaming algorithms in the black-box adversarial model [9,8,37,73,4,48,15,56,6,7,17], in which the adversary can monitor the outputs of the streaming algorithm and choose future stream updates based on these outputs. While useful in a number of applications, there are other settings where the adversary may also have access to the internal state of the algorithm, and this necessitates looking at a stronger adversarial model known as a white-box adversarial stream.

### 1.1 White-box Adversarial Streaming Model

We consider the white-box adversarial streaming model introduced in [2], where a sequence of stream updates is chosen adaptively by an adversary who sees the full internal state of the algorithm at all times, including the parameters and the previous randomness used by the algorithm. The interaction between an adversary and the streaming algorithm is informally depicted as the following game.

*Consider a multi-round, two-player game between **StreamAlg**, the streaming algorithm, and an adversary  $\mathcal{A}$ . Prior to the beginning of the game, fix a query function  $f$ , which asks for a function of an underlying dataset that will be implicitly defined by the stream. In each round:*

- 1.  $\mathcal{A}$  computes an update  $x$  for the stream, which depends on all previous stream updates, all previous internal states, and the randomness used by **StreamAlg**.*
- 2. **StreamAlg** acquires a fresh batch of random bits, uses  $x$  to update its data structures, and (if asked) outputs a response to the query function  $f$ .*
- 3.  $\mathcal{A}$  observes the random bits, the internal state of **StreamAlg**, and the response.*

*The goal of  $\mathcal{A}$  is to make **StreamAlg** output an incorrect query at some round during the game.*

We give a more formal Definition 2.12 in Section 2.3, but give a few motivating applications here.

The white-box adversarial model captures characteristics of many real-world attacks, where an adaptive adversary has access to the entirety of the internal states of the system. In comparison to the oblivious stream model or the black-box adversarial model [8], this model allows for richer adversarial scenarios.

For example, consider a distributed setting where a centralized server collects statistics of a database generated by remote users. The server may send components of its internal state  $S$  to the remote users in order to optimize the total communication over the network. The remote users may use  $S$  in some process that generates downstream data. Thus, future inputs depend on the internal state  $S$  of the algorithm run by the central coordinator. In such settings, the white-box robustness of the algorithms is crucial for optimal selection of query plans [62], online analytical processing [64], data integration [16], and data warehousing [25].

Many persistent data structures provide the ability to quickly access previous versions of information stored in a repository shared across multiple collaborators. The internal persistent data structures used to provide version control may be accessible and thus visible to all users of the repository. These users may then update the persistent data structure in a manner that is not independent of previous states [26,28,47].

Dynamic algorithms often consider an adaptive adversary who generates the updates upon seeing the entire data structure maintained by the algorithm during the execution [18,19,60]. For example, [71] assumes the entire state of the algorithm (including the set of randomness) is available to the adversary after each update, i.e., a white-box model.

Moreover, robust algorithms and adversarial attacks are important topics in machine learning [68,36], with a large body of recent literature focusing on adversarial robustness of machine learning models against white-box attacks [41,55,61,69,24,50,53].

## 1.2 Prior Work in the White-Box Model

We aim to design *white-box adversarially robust* (WAR) streaming algorithms with provable correctness guarantees against such powerful adversaries. In general, it is difficult to design non-trivial WAR algorithms, due to the fact that many widely used space-saving techniques are subject to attacks when all parameters and generated randomness of the system are immediately revealed to the adversary. In the streaming community, both [2] and [27] explicitly study the white-box model for data streams and suggest the use of cryptography in designing WAR algorithms (see also [8] for the use of cryptography for black box robustness). A similar model to the white-box model has also been independently developed in the cryptography community. In property preserving hashing [11] one requires that the hash  $h$  allows, given  $h(\mathbf{x}), h(\mathbf{y})$  for two inputs  $\mathbf{x}, \mathbf{y}$ , to compute a property  $P(\mathbf{x}, \mathbf{y})$  such as distances of  $\mathbf{x}$  and  $\mathbf{y}$  in various norms. The robust versions, similar in spirit to white-box adversarially robust streaming, require that it is computationally hard to find  $\mathbf{x}, \mathbf{y}$  so that from the hashes, one cannot predict the property well. There is some overlap between the properties studied in this model and the problems studied in the white-box model, such as estimating the number of non-zero entries of a vector. We compare the models more below.

A notable drawback of the solutions proposed in both communities is their reliance on either a *random oracle* or a prohibitively long random string. That is, a key component to their algorithms is a random matrix  $\mathbf{A}$  for which it is hard to find a vector  $\mathbf{x}$  with small entries for which  $\mathbf{Ax}$  is sparse or even zero. Unfortunately, it is unknown how to generate  $\mathbf{A}$  pseudorandomly and still maintain white-box adversarial robustness. This is especially problematic in data streams, where the resource measure is the total space complexity, including the randomness, and there is a large body of work on derandomizing such algorithms, see, e.g., [58,42,46,45]. This motivates a central question of this work:

**Question:** Do there exist efficient streaming algorithms in the white-box model without a random oracle?

Another drawback of prior work is the time required of such schemes. The sketching matrices  $\mathbf{A}$  used in [2,27] are random dense matrices, and thus the time to process each stream item, which we refer to as the *update time*, is as large as the number of rows of  $\mathbf{A}$ , which can be large depending on the application. For some problems, such as estimating the number of non-zeros of an underlying vector  $\mathbf{x}$  in a stream up to a constant factor, it was shown that the matrix  $\mathbf{A}$  can be sparse [11]. However, the results of [11] are based on the less studied Sparse Short Vector assumption, require a random oracle or storing the entire matrix  $\mathbf{A}$ , and do not apply to the main problem of sparse recovery that we study. However, by replacing the list-decoding sketch in the followup work [39] with a sparse compressed sensing scheme based on expanders [10], it seems possible to significantly reduce the update time even for the central problem of sparse recovery that we study. Unfortunately, [39] also only seems to work in the random oracle model or require the storage of a very long random string. Thus, achieving fast update time while accounting for the total memory including that to store randomness, is an essential goal.

We now define the central problem we consider in the white-box model, which is the  $k$ -sparse recovery problem as follows:

**Sparse Recovery:** given an input vector  $\mathbf{x}$ , if  $\mathbf{x}$  contains at most  $k$  non-zero entries, recover  $\mathbf{x}$ . Otherwise, report *invalid input*.

A primitive for sparse recovery in the white-box model can be used for tasks such as distinct element estimation, low rank approximation of matrices and tensors, and finding a maximum matching [2,27]. Sparse recovery is also itself a critical data acquisition and processing problem that arises in signal and image processing, machine learning, data networking, and medicine, see [2,27] and the references therein.

We note that a weaker definition of sparse recovery assumes that the input vector  $\mathbf{x}$  is promised to be  $k$ -sparse, and only under this promise recovers  $\mathbf{x}$ . If the input is not  $k$ -sparse, the algorithm can output any vector. This problem can in fact be solved deterministically (and thus also in the white-box model) with low memory (see, e.g., [10], for a version with fast update time, though other simpler schemes based on MDS codes exist). The drawback under this weaker definition

though is that, without knowing the sparsity of the input ahead of time, a user of the algorithm cannot tell whether the output is a correct recovery or garbage. In the latter case, the client might ignorantly use the garbage output in subsequent applications and propagate such an error.

As discussed above, the main issues with previous work on sparse recovery in the white-box setting [2,27] or related property preserving hashing works [11,39] is that they rely on a random oracle, or do not apply to the sparse recovery problem itself, or do not have fast update time.

### 1.3 Our Results

We let  $n$  be the dimension of the underlying input vector, and for notational convenience, we assume the stream length is at most a polynomial in  $n$ .

We also assume for notational convenience that the entries of the vector are integers bounded in absolute value by  $\text{poly}(n)$  at all times during the stream. We use  $\tilde{O}(f)$  to denote  $f \cdot \text{poly}(\log n)$ . We denote the security parameter of our cryptographic assumption as  $\lambda$ . Our robustness guarantees hold against adversaries that run in time polynomial in  $n$ . In our applications, the efficiency metrics will depend on the security parameter  $\lambda$  used for the cryptographic components. We can set  $\lambda = n^c$  for an arbitrary small constant  $c > 0$ . Further, assuming subexponential time assumptions we could also set  $\lambda$  to be a large enough polylogarithmic function of  $n$ . Let  $\ell(\lambda)$  denote a fixed polynomial in  $\lambda$  that represents the maximum run-time of various cryptographic algorithms involved. By setting  $\lambda = n^c$  appropriately for a small enough constant  $c$  we can always set  $\ell(\lambda)$  to be a small enough polynomial in  $n$  (or a polylogarithmic function of  $n$  assuming subexponential time assumptions). Our results are described in terms of  $\ell(\lambda)$ .

Our main result is informally stated as follows:

**Theorem 1.1 (Informal).** *Assuming the hardness of the Learning with Errors (LWE) problem, there is a WAR streaming algorithm for  $k$ -sparse recovery, which:*

- takes  $\tilde{O}(k + \ell(\lambda))$  bits of space,
- has  $\tilde{O}(1 + \ell(\lambda))$  update time, and
- has  $\tilde{O}(k \cdot (k^c + \ell(\lambda)))$  report time for an arbitrarily small constant  $c > 0$ .

This matches (up to a poly-logarithmic factor) the optimal space and update time complexity for the weaker  $k$ -sparse recovery problem under the oblivious streaming setting. On the other hand, assuming polynomial security of LWE, there is a multiplicative overhead of  $n^\epsilon$  for arbitrary constant  $\epsilon > 0$ . This matches the complexity of the previous work [27] that was proven to be secure in the random oracle model.

We also study the sparse recovery problem in a white-box distributed model, which is a computational framework that is closely related to the data stream model. See Section A.1 for a formal discussion of the model, which involves multiple servers who communicate with a so-called coordinator. A WAR streaming

algorithm can be naïvely converted to a WAR distributed algorithm under certain conditions, yet the resulting processing time of the servers is  $\mathcal{O}(n)$  times that of the streaming update time. We show that we can achieve a near-optimal processing time in a white-box distributed model assuming the *polynomial* hardness of Ring-LWE. This is in contrast to the suboptimal time implied by the naïve conversion above and we also do not require a subexponential hardness assumption. This is due to the fact that Ring-LWE supports SIMD (Single-Instruction Multiple-Data) style operations. We informally state our result as follows:

**Theorem 1.2 (Informal).** *Assuming the hardness of Ring-LWE, there exists a WAR distributed protocol for  $k$ -sparse recovery, which:*

- *uses  $\tilde{\mathcal{O}}(k + \ell(\lambda))$  bits of communication between each server and the coordinator,*
- *has  $\tilde{\mathcal{O}}(n)$  processing time on each server, and*
- *has the coordinator spend  $\tilde{\mathcal{O}}(\max(n, k^{1+c}))$  time to output the solution given the messages from the servers, for an arbitrarily small constant  $c > 0$ .*

We summarize the complexities of our algorithms in Table 1, as compared to the best-known upper bounds for these problems in the white-box adversarial model.

**Table 1.** A summary of the bit complexities and runtime of our  $k$ -sparse recovery algorithm, as compared to the best-known upper bounds for these problems in the white-box adversarial model. The TIME columns refer to the update time in the streaming model and processing time in the distributed model.  $c$  denotes an arbitrarily small positive constant.

	MODEL	ASSUMPTION	SPACE	TIME	RAND. ORACLE?
<b>Prev.</b>	STREAMING	SUBEXP SIS	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(k)$	✓
	DISTRIBUTED	SUBEXP SIS	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(nk)$	✓
<b>Ours</b>	STREAMING	SUBEXP LWE	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(1)$	✗
	STREAMING	POLY LWE	$\tilde{\mathcal{O}}(n^c + k)$	$\tilde{\mathcal{O}}(n^c)$	✗
	DISTRIBUTED	SUBEXP LWE	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(n)$	✗
	DISTRIBUTED	POLY LWE	$\tilde{\mathcal{O}}(n^c + k)$	$\tilde{\mathcal{O}}(n^{1+c})$	✗
	DISTRIBUTED	POLY RING-LWE	$\tilde{\mathcal{O}}(n^c + k)$	$\tilde{\mathcal{O}}(n)$	✗

Similar to [27], generalizing our construction produces low-rank matrix and tensor recovery algorithms, again with the crucial property that they do not assume a random oracle or store a large random seed.

Consider a data stream updating the entries of an underlying matrix or tensor, which can be represented as its vectorization in a stream. Let  $n \geq m$  be the dimensions of an underlying input matrix, and let  $n$  be the dimension

parameter of an underlying data tensor in  $\mathbb{R}^{n^d}$  for some  $d \in \mathcal{O}(1)$  (we assume all dimensions of the tensor are the same only for ease of notation). We have the following results in the streaming model:

**Theorem 1.3 (Informal).** *Assuming the hardness of the Learning with Errors (LWE) problem, there is a WAR streaming algorithm for rank- $k$  matrix recovery, which takes  $\tilde{\mathcal{O}}(nk + \ell(\lambda))$  bits of space. There is also a WAR streaming algorithm for rank- $k$  tensor recovery, which takes  $\tilde{\mathcal{O}}(nk^{\lceil \log d \rceil} + \ell(\lambda))$  bits of space.*

We remark that the above matrix recovery algorithm performs  $\tilde{\mathcal{O}}(nk)$  measurements using  $n$ -sparse matrices (plus a relatively small FHE hash) which can be converted into a distributed protocol with  $\tilde{\mathcal{O}}(n^2k)$  processing time. This improves upon the matrix recovery algorithm in the previous work [27] which, if implemented in a distributed setting, would require  $\tilde{\mathcal{O}}(n^3k)$  processing time at each server.

*Additional Applications.* Our  $k$ -sparse recovery algorithm can be used to construct an  $L_0$ -norm estimation algorithm following the same recipe in [27]. In short, by running our  $k$ -sparse recovery algorithm with the sparsity parameter  $k$  set to  $n^{1-w}$  for some constant  $w$ , we can get an  $n^w$ -approximation to the  $L_0$  norm of a vector. In comparison to [27], our construction again removes the need for a random oracle or the need to store a lot of randomness.

In addition, our construction implies a family of Homomorphic Collision-Resistance hash functions defined in [39], with concrete efficiency, security against white-box adversaries (i.e., adversaries with access to the randomness of sampling the hash function), and especially short hash function descriptions, which may be of independent interest.

#### 1.4 Technical Overview

We leverage structural properties of lattice-based assumption to address all drawbacks listed above. In particular, departing from the previous construction of [27], we work with a suitable FHE scheme to achieve the desired property. We start with an intuitive explanation of the previous construction in [27].

*Previous Construction* One general framework for solving the sparse recovery problem is to run a weak recovery algorithm as a black-box subroutine, which potentially provides garbage output, in addition to an additional “tester” scheme to verify whether the recovery output matches the original input (i.e., is a correct recovery) or is garbage. In the oblivious streaming setting, such a “tester” can be implemented using a space and time-efficient randomized sampler [23]. However, this approach is not WAR as the randomness of the sampler is revealed and thus not independent of the stream.

Instead, [27] uses a random matrix  $\mathbf{A}$  to implement a WAR tester, which is assumed to be heuristically compressed by an idealized hash function. During the stream, [27] maintains a hash  $h = \mathbf{A}\mathbf{x}$ , where  $\mathbf{x}$  is the data vector implicitly

defined by the stream. This is easy to maintain since  $\mathbf{A}$  is linear and the stream undergoes additive updates of the form  $\mathbf{x}_i \leftarrow \mathbf{x}_i + \Delta$  to the coordinates  $\mathbf{x}_i$  of  $\mathbf{x}$ . In the testing stage, it acknowledges an output  $\mathbf{x}'$  to be a correct recovery if and only if  $\mathbf{x}'$  is  $k$ -sparse and  $\mathbf{A}\mathbf{x}'$  equals the hash  $h$ . The idea is that if there exists a  $k$ -sparse collision  $\mathbf{x}'$  satisfying  $\mathbf{A}(\mathbf{x} - \mathbf{x}') = 0$  and  $\mathbf{x}' \neq \mathbf{x}$ , then an adversary can brute-force over all  $k$ -sparse vectors to find such  $\mathbf{x} - \mathbf{x}'$ , which is a solution to the SIS problem. To spawn  $n$  matrix columns and support the brute-force reduction, the security parameter (i.e., the number of rows of  $\mathbf{A}$ ) needs to be in  $\omega(k \log n)$ , assuming the exponential security of SIS.

*Efficient Reduction.* The number of matrix rows in the previous construction depends on the sparsity parameter  $k$  due to a brute-forcing step in the reduction. To grant enough time for iterating through all  $k$ -sparse vectors (with  $\text{poly}(n)$ -bounded entries), the security parameter has to be greater than  $k \log n$ , which results in  $\tilde{O}(k)$  time to process every single update. Our first idea is to construct a poly-time reduction in order to remove such a dependency on  $k$ . This allows us to achieve better space and time efficiency, in addition to basing the construction on a polynomial hardness assumption.

The observation is that the non-existence of a sparse collision is an overly strong condition. Instead, for our purpose, we only need to ensure that such a collision, if it exists, cannot be found by the poly-time weak recovery scheme that the algorithm runs. Thus, in the reduction, instead of brute-forcing over all  $k$ -sparse vectors, we let the adversary run the same poly-time weak recovery scheme that the algorithm does. If a collision  $\mathbf{x}'$  satisfying  $\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{x}'$  is output by the scheme, then the adversary can still break the security assumption.

*Reducing the Amount of Randomness.* The rest of our effort will be attributed to reducing the size of the state of the algorithm. Rather than using a truly random giant matrix to hash the data as in the previous construction, we would like to design a short digest (ideally linear-sized in the security parameter) and generate the columns of the hash matrix on the fly. This naturally yields correlated columns as opposed to the truly random columns in an SIS matrix. We aim to show that in this setting, the chance of finding a collision is still negligible.

*Streaming Friendly Collision-Resistant Hash* Our main insight from cryptography to address the issue above is that we could construct a hash function family  $\mathcal{H}$ , with special properties to help solve our streaming problem:

- The hash key  $h \leftarrow \mathcal{H}$  is small in size, ideally polynomial in the security parameter (or even polylogarithmic assuming subexponential time assumptions), in particular independent of the dimension  $n$  of the stream vector.
- The hash should be computable in a streaming fashion and it should be updatable where each update takes polynomial time in  $\lambda$  (or even polylogarithmic time assuming subexponential time assumptions) independent of  $n$  or  $k$  or other parameters.



- The hash key should be pseudorandom. In other words, in any real world implementation in the streaming setting the model does not allow us to use any structured randomness (such as with hidden planted secrets). In the proof however, we could introduce these secrets which would be important for our construction.
- The hash should satisfy collision-resistance.

*Constructing Streaming Friendly Hashing from FHE.* How could we build such a hash function? The SIS based hash function does have a fast update time. On the other hand, we could take an arbitrary collision resistant function that compresses the number of bits by a factor of  $1/2$  and use it to perform a Merkle style hashing to hash strings of arbitrary length. Such a hash will have  $\text{poly}(\lambda)$  size hashing keys independent of the stream length. Unfortunately, the first idea suffers from a large hashing key size, which could be shortened in the random order (RO) model by using RO to store the SIS matrix. Moreover, the second proposal will not support positive and negative stream updates to the coordinates of the underlying vector.

Our main idea is a modular construction that leverages a suitably chosen FHE scheme satisfying mild structural properties (that are satisfied currently by most schemes based on LWE and Ring-LWE) to build such a hash function. Our ideas are somewhat inspired from [14], which constructed ABE with shorter public keys to encrypt with respect to attributes of unbounded length.

We directly explain the idea in the context of our construction.

In our hash function, the hash key consists of  $L = O(\log n)$  ciphertexts  $\text{ct}_1, \dots, \text{ct}_L$ . We assume that our FHE scheme has pseudorandom ciphertexts and keys. In our actual scheme, these ciphertexts will be generated as random strings. It is only in the proof that they will correspond to encryptions of carefully chosen values.

The idea is that from these  $L$  ciphertexts, we will derive  $n$  evaluated ciphertexts  $\widehat{\text{ct}}_j$  for  $j \in [n]$  by evaluating  $L$  ciphertexts on simple functions  $C_1, \dots, C_n$  (described later). These  $n$  ciphertexts will correspond to columns of the SIS type hashing. The hash of a vector  $\mathbf{x}$  is computed as the linear combination  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i$  where  $\mathbf{x}_i$  is the  $i^{\text{th}}$  coordinate. Here the sum is done over the field on which the FHE ciphertexts live.

The circuits  $C_1, \dots, C_n$  are chosen keeping the security proof in mind. The ciphertexts  $\text{ct}_1, \dots, \text{ct}_L$  are “thought” to encrypt a value  $m \in \{0, 1\}^L$  that will be chosen in the proof, and each  $C_i(\cdot)$  on input  $m$  outputs 1 if  $i = m$  and 0 otherwise. Observe that these circuits are polynomial in  $L = O(\log n)$  sized and therefore can be evaluated on  $\text{ct}_1, \dots, \text{ct}_L$  in polynomial in  $(\lambda, \log n)$  time guaranteeing fast updates.

How do we prove security? The idea is that if an adversary breaks collision-resistance it will produce two vectors  $\mathbf{x} \neq \mathbf{x}'$ , so that the corresponding hashes  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i = \sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}'_i$ . Suppose one could know exactly one coordinate  $v \in [n]$  so that  $\mathbf{x}_v \neq \mathbf{x}'_v$ . Then we could set  $m = v$ . In this case  $\widehat{\text{ct}}_v$  encrypts 1 and the other ciphertexts encrypt 0. In this case,  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i = \sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}'_i$

would lead to a contradiction if our FHE scheme has a special property (that is satisfied by common FHE schemes). The property is that if  $\widehat{\text{ct}}_i$  encrypts bit  $\mu_i \in \{0, 1\}$  and  $\mathbf{x}_i$  are small norm then there is a decryption algorithm that uniquely binds  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i$  to  $\sum_{i \in [n]} \mu_i \mathbf{x}_i$ . If this holds, then the equality  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i = \sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}'_i$  would imply  $\mathbf{x}_v = \mathbf{x}'_v$  which is a contradiction.

This leads us to two remaining problems:

1. Given that we do not know  $\mathbf{x}$  and  $\mathbf{x}'$  ahead of time, how do we decide which index  $m$  to encrypt such that it will help identify whether  $\mathbf{x} = \mathbf{x}'$ ?
2. Moreover, a WAR algorithm cannot actually encrypt FHE ciphertexts, since all randomness and parameters used by the algorithm will be exposed to the adversary. Instead, we want our digest to be unstructured and truly random.

*FHE with a Pseudorandom Property.* We solve both issues above by considering a pseudorandom property of the FHE scheme and show that a random guess of the index  $m$  actually suffices. We assume that the distributions of public keys and ciphertexts of the FHE scheme are indistinguishable, respectively, from some truly random distributions from the perspective of the adversary.

We sample truly random digests in our construction, and perform a random guess of  $m \in [n]$  in the proof. If an adversary finds a collision pair  $\mathbf{x} \neq \mathbf{x}'$  with probability  $\gamma$ , then with probability  $p \geq \gamma/n$ , our guess pins down a coordinate of disagreement  $\mathbf{x}_m \neq \mathbf{x}'_m$ . In the proof, we then switch the ciphertexts to encrypt  $m$ . The chance of guessing the disagreement is still  $p$ , as otherwise one can distinguish the pseudorandom ciphertext. However, as argued earlier,  $\mathbf{x}_m \neq \mathbf{x}'_m$  roughly implies  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i \neq \sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}'_i$ , and thus two distinct vectors cannot form a hash collision that fools our “tester”.

In our complexity analysis, we instantiate the FHE scheme using GSW, which satisfies the additional properties required by our construction. We remark that in addition to GSW [33], most current FHE schemes based on LWE and Ring-LWE also exhibit these properties, and thus they are not an artifact obtained from a particular FHE scheme.

*Ring-LWE for the Distributed Setting.* In the distributed model, each remote server receives as input a length- $n$  vector. The goal is to design a protocol such that each server can efficiently compute a short summary of their partition of data. Adapting the above streaming algorithm to this setting, each server needs to evaluate  $n$  ciphertexts in order to compute a linear combination  $\sum_{i \in [n]} \widehat{\text{ct}}_i \mathbf{x}_i$ .

Using a polynomially secure LWE-based FHE scheme, it is unclear how to speed up such  $n$ -many individual evaluations as each computation takes  $\text{poly}(\lambda)$  operations, resulting in total complexity of  $n \text{poly}(\lambda)$  where  $\lambda$  is the security parameter. However, appealing to the SIMD property of Ring-LWE and packing the data as ring elements, we can represent the input vectors as partitioned segments of length proportional to the security parameter (rather than  $n$  individual coordinates), and efficiently operate on them segment-by-segment using the Fast Fourier Transform. This improves the processing time of each server to be almost linear in the dimension of the data. We refer the reader to Section A.1 for more details.

*Summary.* Leveraging FHE or lattice-based assumptions in algorithm design are relatively new techniques. We think it would be good to explore their applications more broadly in robust algorithms, beyond the recovery problems or even the adversarial settings considered here.

## 2 Preliminaries

*Basic Notations.* Let  $\mathbb{1}(p)$  for some proposition  $p$  be an indicator variable, which equals 1 if  $p$  is true and 0 otherwise. Define  $[a, b] := \{a, a+1, \dots, b\}$ ,  $[n] := [1, n]$ . For a finite set  $S$ , we write  $x \xleftarrow{\$} S$  to denote uniformly sampling  $x$  from  $S$ . Our logarithms are in base 2. For  $n \in \mathbb{R}$ , we use  $\lceil n \rceil$  to denote rounding  $n$  up to the nearest integer. For a vector  $\mathbf{x}$ ,  $\|\mathbf{x}\|_0$  denotes the  $\ell_0$  norm of  $\mathbf{x}$ , which is the number of its non-zero entries. We abbreviate PPT for probabilistic polynomial-time. A function  $\text{negl} : \mathbb{N} \rightarrow \mathbb{R}$  is negligible if for every constant  $c > 0$  there exists an integer  $N_c$  such that  $\text{negl}(x) < x^{-c}$  for all  $x > N_c$ .

### 2.1 Lattice Preliminaries

*General definitions.* A lattice  $\mathcal{L}$  is a discrete subgroup of  $\mathbb{R}^g$ , or equivalently the set  $\mathcal{L}(\mathbf{b}_1, \dots, \mathbf{b}_m) = \{\sum_{i=1}^g x_i \mathbf{b}_i : x_i \in \mathbb{Z}\}$  of all integer combinations of  $g$  vectors  $\mathbf{b}_1, \dots, \mathbf{b}_m \in \mathbb{R}^g$  that are linearly independent. Such  $\mathbf{b}_i$ 's form a *basis* of  $\mathcal{L}$ . The lattice  $\mathcal{L}$  is said to be *full-rank* if  $g = m$ . We denote by  $\lambda_1(\mathcal{L})$  the so-called first minimum of  $\mathcal{L}$ , defined to be the length of a shortest non-zero vector of  $\mathcal{L}$ .

*Discrete Gaussian and Related Distributions* For any  $s > 0$ , define  $\rho_s(\mathbf{x}) = \exp(-\pi\|\mathbf{x}\|^2/s^2)$  for all  $\mathbf{x} \in \mathbb{R}^g$ . We write  $\rho$  for  $\rho_1$ . For a discrete set  $S$ , we extend  $\rho$  to sets by  $\rho_s(S) = \sum_{\mathbf{x} \in S} \rho_s(\mathbf{x})$ . Given a lattice  $\mathcal{L}$ , the *discrete Gaussian*  $\mathcal{G}_{\mathcal{L},s}$  is the distribution over  $\mathcal{L}$  such that the probability of a vector  $\mathbf{y} \in \mathcal{L}$  is proportional to  $\rho_s(\mathbf{y})$ :  $\Pr_{X \leftarrow \mathcal{G}_{\mathcal{L},s}}[X = \mathbf{y}] = \frac{\rho_s(\mathbf{y})}{\rho_s(\mathcal{L})}$ .

The discrete Gaussian is not bounded but with overwhelming probability a sample from  $\mathcal{G}_{\mathcal{L},s}$  is bounded by  $s\sqrt{g}$  in  $\ell_2$  norm.

*Lattice problems.* We consider worst-case and average-case problems over lattices. We start by defining the gap version of the shortest vector problem.

**Definition 2.1.** For any function  $\gamma = \gamma(g) \geq 1$ , the decision problem  $\text{GapSVP}_\gamma$  (Gap Shortest Vector Problem) is defined as follows: the input is a basis  $\mathbf{B}$  for a lattice  $\mathcal{L} \subset \mathbb{R}^g$  and a real number  $d > 0$  that satisfy the promise that  $\lambda_1(\mathcal{L}(\mathbf{B})) \leq d$ , or  $\lambda_1(\mathcal{L}(\mathbf{B})) \geq \gamma d$ . The goal is to decide whether  $\lambda_1(\mathcal{L}(\mathbf{B})) \leq d$ , or  $\lambda_1(\mathcal{L}(\mathbf{B})) \geq \gamma d$ .

This problem is known to be NP hard for  $\gamma = g^{o(1)}$  [49,38], and is known to be inside  $\text{NP} \cap \text{coNP}$  for  $\gamma = O(\sqrt{g})$  [1,35]. The problem is believed to be intractable by polynomial time adversaries even for  $\gamma = 2^{g^\epsilon}$  for  $0 < \epsilon < 1$ .

Next, we define the Learning with Errors problem LWE.

**Definition 2.2 (Learning with Errors).** For any integer dimension  $g \in \mathbb{N}$ , modulus  $q \in \mathbb{N}$ , sample complexity  $h \geq \Omega(g \cdot \log q)$  and a parameter  $\sigma \in \mathbb{R}^{\geq 0}$ . Define an error distribution as the discrete Gaussian  $\chi_\sigma^h := \mathcal{G}_{\sigma, \mathbb{Z}^h}$ . We define  $\text{LWE}_{g,h,q,\chi_\sigma^h}$  as follows: Sample a matrix  $\mathbf{A} \leftarrow \mathbb{Z}_q^{g \times h}$  and a secret  $\mathbf{s} \in \mathbb{Z}_q^{1 \times g}$ . Sample an error vector  $\mathbf{e} \leftarrow \chi_\sigma^h$ . The problem is to distinguish  $(\mathbf{A}, \mathbf{s}\mathbf{A} + \mathbf{e} \bmod q)$  from  $(\mathbf{A}, \mathbf{r})$  where  $\mathbf{r}$  is a random in  $\mathbb{Z}_q^{1 \times h}$ .

Under quantum and classical reductions it can be shown that  $\text{LWE}_{g,h,q,\chi_\sigma}$  is harder to solve than  $\text{Gap-SVP}_\gamma$  with  $\gamma = g \frac{q}{\sigma}$ .

We also consider the Ring-LWE problem, which is an average-case lattice problems over ideal lattices.

**Definition 2.3 (Ring-LWE).** Let  $q \geq 1$  be a prime number, and let  $g$  be a power of 2. Let  $g, \sigma$  be positive integers. The  $\text{Ring-LWE}_{g,h,q,\chi_\sigma}$  problem is defined as follows: Consider  $R = \frac{\mathbb{Z}_q[x]}{x^g + 1}$ . Sample at random  $\mathbf{a}_1, \dots, \mathbf{a}_h \leftarrow R$  and a secret polynomial  $\mathbf{s} \leftarrow R$ . Sample  $h$  error polynomials  $\mathbf{e}_1, \dots, \mathbf{e}_h$  such that each  $\mathbf{e}_i$  is a  $g$ -dimensional vector  $\mathbf{e}_i \leftarrow \chi_\sigma = \mathcal{G}_{\sigma, \mathbb{Z}^g}$ . The problem is to distinguish the tuple  $(\mathbf{a}_1, \dots, \mathbf{a}_h, \{\mathbf{s}\mathbf{a}_i + \mathbf{e}_i\}_{i \in [h]})$  from  $(\mathbf{a}_1, \dots, \mathbf{a}_h, \{\mathbf{r}_i\}_{i \in [h]})$  where  $\mathbf{r}_1, \dots, \mathbf{r}_h$  are random polynomials in  $R$ .

The Ring-LWE problem is as hard as worst-case problems over ideal lattices [54]. Ring LWE is sometimes advantageous to use over LWE because of shorter parameters. One ring LWE element is of dimension  $g$ , and requires only one ring element as a coefficient, as opposed to  $g$  vectors in the case of LWE. Moreover, ring multiplication can be sped up using number-theoretic transforms [65].

## 2.2 Fully Homomorphic Encryption

We review the definition of fully homomorphic encryption.

*Notation.* We denote the security parameter of a cryptographic scheme by  $\lambda$ . We write  $\text{PPT}_\lambda$  to denote probabilistic poly-time in  $\lambda$ . We say two ensembles of distributions  $\{\mathcal{D}_\lambda\}_{\lambda \in \mathbb{N}}$  and  $\{\mathcal{D}'_\lambda\}_{\lambda \in \mathbb{N}}$  are computationally indistinguishable, denoted by  $\approx_c$ , if for any non-uniform PPT adversary  $\mathcal{A}$  there exists a negligible function  $\text{negl}$  such that  $\mathcal{A}$  can distinguish between the two distributions with probability at most  $\text{negl}(\lambda)$ .

*Remark 2.4.* We also define a  $(\mathcal{T}, \epsilon)$ -indistinguishability, where  $\mathcal{T}$  and  $\epsilon$  are two functions in  $\lambda$ . We say that two distribution ensembles  $\{\mathcal{D}_\lambda\}_{\lambda \in \mathbb{N}}$  and  $\{\mathcal{D}'_\lambda\}_{\lambda \in \mathbb{N}}$  are  $(\mathcal{T}, \epsilon)$ -indistinguishable, denoted by  $\approx_{(\mathcal{T}, \epsilon)}$ , if there exists an integer  $\lambda_0$  such that for all  $\lambda > \lambda_0$ , for every adversary  $\mathcal{A}$ ,  $\mathcal{A}$  can distinguish between the two distributions with probability at most  $\epsilon(\lambda)$ .

**Definition 2.5 (FHE).** A fully homomorphic encryption FHE scheme is a tuple of  $\text{PPT}_\lambda$  algorithms  $\text{FHE} = (\text{FHE.Setup}, \text{FHE.Enc}, \text{FHE.Eval}, \text{FHE.Dec})$  with the following properties:

- $\text{FHE.Setup}(1^\lambda) \rightarrow (\text{pk}, \text{sk})$  : On input a security parameter  $\lambda$ , the setup algorithm outputs a key pair  $(\text{pk}, \text{sk})$ .
- $\text{FHE.Enc}(\text{pk}, \mu) \rightarrow \text{ct}$  : On input a public key  $\text{pk}$  and a message  $\mu \in \{0, 1\}$ , the encryption algorithm outputs a ciphertext  $\text{ct}$ .
- $\text{FHE.Eval}(\text{pk}, C, \text{ct}_1, \dots, \text{ct}_\ell) \rightarrow \hat{\text{ct}}$  : On input a public key  $\text{pk}$ , a circuit  $C : \{0, 1\}^\ell \rightarrow \{0, 1\}$  of depth at most  $\text{poly}(\lambda)$ , and a tuple of ciphertexts  $\text{ct}_1, \dots, \text{ct}_\ell$ , the evaluation algorithm outputs an evaluated ciphertext  $\hat{\text{ct}}$ .
- $\text{FHE.Dec}(\text{pk}, \text{sk}, \hat{\text{ct}}) \rightarrow \hat{\mu}$  : On input a public key  $\text{pk}$ , a secret key  $\text{sk}$ , and a ciphertext  $\hat{\text{ct}}$ , the decryption algorithm outputs a message  $\hat{\mu} \in \{0, 1, \perp\}$ .

We also require an FHE scheme to satisfy compactness, correctness, and security properties as follows:

**Definition 2.6 (FHE Compactness).** We say that a FHE scheme is compact if there exists a polynomial  $\text{poly}(\cdot)$  such that for all security parameter  $\lambda$ , circuit  $C : \{0, 1\}^\ell \rightarrow \{0, 1\}$  of depth at most  $\text{poly}(\lambda)$ , and  $\mu_i \in \{0, 1\}$  for  $i \in [\ell]$ , the following holds: Given  $(\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)$ ,  $\text{ct}_i \leftarrow \text{FHE.Enc}(\text{pk}, \mu_i)$  for  $i \in [\ell]$ ,  $\hat{\text{ct}} \leftarrow \text{FHE.Eval}(\text{pk}, C, \text{ct}_1, \dots, \text{ct}_\ell)$ , we always have  $|\hat{\text{ct}}| \leq \text{poly}(\lambda)$ .

**Definition 2.7 (FHE Correctness).** We say that a FHE scheme is correct if for all security parameter  $\lambda$ , circuit  $C : \{0, 1\}^\ell \rightarrow \{0, 1\}$  of depth at most  $\text{poly}(\lambda)$ , and  $\mu_i \in \{0, 1\}$  for  $i \in [\ell]$ , the following holds: Given  $(\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)$ ,  $\text{ct}_i \leftarrow \text{FHE.Enc}(\text{pk}, \mu_i)$  for  $i \in [\ell]$ ,  $\hat{\text{ct}} \leftarrow \text{FHE.Eval}(\text{pk}, C, \text{ct}_1, \dots, \text{ct}_\ell)$ ,

$$\Pr[\text{FHE.Dec}(\text{pk}, \text{sk}, \hat{\text{ct}}) = f(\mu_1, \dots, \mu_\ell)] = 1.$$

**Definition 2.8 (FHE Security).** We say that a FHE scheme satisfies security if for all security parameter  $\lambda$ , the following holds: For any  $\text{PPT}_\lambda$  adversary  $\mathcal{A}$ , there exists a negligible function  $\text{negl}$  such that

$$\text{adv}_{\mathcal{A}, \text{FHE}} := 2 \cdot |1/2 - \Pr[\text{Expt}_{\mathcal{A}, \text{FHE}}(1^\lambda) = 1]| < \text{negl}(\lambda)$$

, where the experiment  $\text{Expt}_{\mathcal{A}, \text{FHE}}$  is defined as follows:

$\text{Expt}_{\mathcal{A}, \text{FHE}}(1^\lambda)$  :

1. On input the security parameter  $1^\lambda$ , the challenger runs  $(\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)$ , and  $\text{ct} \leftarrow \text{FHE.Enc}(\text{pk}, b)$  for  $b \xleftarrow{\$} \{0, 1\}$ . It provides  $(\text{pk}, \text{ct})$  to  $\mathcal{A}$ .
2.  $\mathcal{A}$  outputs a guess  $b'$ . The experiment outputs 1 if and only if  $b = b'$ .

*Remark 2.9.* The security definition above is in terms of  $\text{PPT}_\lambda$  adversaries and  $\text{negl}$  advantage, but it can be lifted to a  $(\mathcal{T}, \epsilon)$  setting similar to Remark 2.4 as follows: We say that a FHE scheme is  $(\mathcal{T}, \epsilon)$ -secure if there exists an integer  $\lambda_0$ , such that for all  $\lambda > \lambda_0$ , for any non-uniform probabilistic  $\mathcal{T}(\lambda)$  time adversary, the advantage  $\text{adv}_{\mathcal{A}, \text{FHE}}$  of the above experiment is at most  $\epsilon(\lambda)$ .

*Additional properties.* Our construction of robust streaming algorithms will be based on an FHE scheme satisfying a few additional properties. We refer to this type of FHE scheme as a *pseudorandom FHE (PFHE) scheme*.

**Definition 2.10 (Pseudorandom FHE (PFHE)).** *Let FHE be a fully homomorphic encryption scheme (Definition 2.5). We call FHE a  $(\mathcal{T}, \epsilon)$ -pseudorandom FHE scheme if for all security parameter  $\lambda$ , it satisfies the following additional properties:*

- FHE.Eval is deterministic.
- (Pseudorandom pk). *There exists an explicit distribution  $\mathcal{D}_{\text{pk}(\lambda)}$ , such that the following two distributions are  $(\mathcal{T}, \epsilon)$ -indistinguishable:*

$$\{\text{pk} : (\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)\} \approx_{(\mathcal{T}, \epsilon)} \mathcal{D}_{\text{pk}(\lambda)}$$

- (Pseudorandom ct). *There exists an explicitly distribution  $\mathcal{D}_{\text{ct}(\lambda)}$ , such that the following two distributions are  $(\mathcal{T}, \epsilon)$ -indistinguishable:*

$$\left\{ \begin{array}{l} (\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda) \\ \text{ct} : \quad \mu \in \{0, 1\} \\ \quad \text{ct} \leftarrow \text{FHE.Enc}(\text{pk}, \mu) \end{array} \right\} \approx_{(\mathcal{T}, \epsilon)} \mathcal{D}_{\text{ct}(\lambda)}$$

- ( $\mathbb{Z}$  Linear Homomorphism). *Given  $(\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)$ , the set of all ciphertexts is a subset of a vector space over  $\mathbb{Z}_q$  for some  $q = q(\lambda)$ . Let  $\mathcal{M} = \sum_{i \in [k]} x_i \cdot \text{ct}_i \bmod q$  be an arbitrary linear combination of ciphertexts. Denote  $\mu_i \leftarrow \text{FHE.Dec}(\text{pk}, \text{sk}, \text{ct}_i)$ . Suppose  $\mathcal{M}$  satisfies:*

1.  $x_i \in \mathbb{Z}$  for all  $i \in [k]$
2.  $k \cdot \max_{i \in [k]} x_i \ll q$

*then we have the followings:  $|\mathcal{M}| \leq \text{poly}(\lambda)$ . And there exists a  $\text{poly}(\lambda)$ -time decoding algorithm LinearDec, such that  $\Pr[\text{LinearDec}(\text{pk}, \text{sk}, \mathcal{M}) = \sum_{i \in [k]} x_i \cdot \mu_i] = 1$ .*

*As in a standard FHE scheme, we also require that a PFHE scheme satisfies compactness (Definition 2.6), correctness (Definition 2.7), and security (Definition 2.8).*

There exist FHE constructions based on LWE and Ring-LWE [33,12,13] assumptions. We will focus on the GSW scheme in our streaming construction and the BV scheme in the distributed setting. Both of them satisfy the above definition (or a slight variant) of pseudorandom FHE.

### 2.3 Streaming Algorithms

We review the streaming model, the structure of a streaming algorithm, and the notion of white-box adversarial robustness.

*Notation.* We denote the input dimension parameter of an algorithm by  $n$ . We write  $\text{PPT}_\lambda$  to denote probabilistic poly-time in  $\lambda$ .

*Streaming Model.* The streaming model captures key resource requirements of algorithms for database, machine learning, and network tasks, where the size of the data is significantly larger than the available storage. This model was formalized by Alon, Matias, and Szegedy [5] as a vector undergoing additive *updates* to its coordinates. In this model, the input data takes the form of a sequence of updates. These updates assume the existence of a fixed-length underlying vector which represents the current dataset, and dynamically modify this vector one coordinate at a time.

Given a dimension parameter  $n \in \mathbb{N}$ , the model behaves as follows:

- At the beginning of the stream, an underlying vector  $\mathbf{x}$  is initialized to  $0^n$ .
- The input stream is a sequence of  $T \leq \text{poly}(n)$  additive updates  $(x_1, \dots, x_T)$ . Each update  $x_t$  for  $t \in [T]$  is a tuple  $x_t := (i_t, \Delta_t)$ , where  $i_t \in [n]$  denotes a coordinate of  $\mathbf{x}$  and  $\Delta_t$  is an arbitrary value.  
 $x_t := (i_t, \Delta_t)$  is interpreted as updating  $\mathbf{x}_{i_t} \leftarrow \mathbf{x}_{i_t} + \Delta_t$ . Thus at any time  $t' \in [T]$  during the stream,  $\mathbf{x}$  is the accumulation of all past updates, with its  $i$ -th coordinate for  $i \in [n]$ :  $\mathbf{x}_i = \sum_{t \leq t'} \mathbb{1}(i_t = i) \cdot \Delta_t$ .

A streaming algorithm receives the stream of updates and uses limited memory to compute a function of  $\mathbf{x}$ . We stress that during the stream, the underlying vector  $\mathbf{x}$  is not explicitly available to the algorithm, but is instead implicitly defined by the sequence of past updates.

In this work, we consider a bounded integer stream, which assumes that throughout the stream,  $\mathbf{x}$  is promised to be in the range  $\{-N, -N+1, \dots, N-1, N\}^n$  for some integer  $N \in \text{poly}(n)$ .

*Streaming Algorithm.* A streaming algorithm maintains some internal state during the stream, which evolves upon receiving new updates. It takes one pass over all updates and then uses its maintained state to compute a function of the underlying dataset  $\mathbf{x}$ . In the streaming model, the dimension parameter  $n$  (the length of  $\mathbf{x}$ ) is typically large. Therefore, a streaming algorithm usually optimizes the size of the state such that it takes sublinear bits of space in  $n$ .

We define the general framework for streaming algorithms as follows:

**Definition 2.11 (Streaming Algorithm).** *Given a query function ensemble  $\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow X\}_{n \in \mathbb{N}}$  for some domain  $X$ , a streaming algorithm  $\text{StreamAlg}$  that computes  $\mathcal{F}$  contains a tuple of operations  $\text{StreamAlg} = (\text{StreamAlg.Setup}, \text{StreamAlg.Update}, \text{StreamAlg.Report})$ , with the following properties:*

- $\text{StreamAlg.Setup}(n)$  : On input a dimension parameter  $n$ , the setup operation initializes an internal data structure  $\mathcal{DS}$  within the memory of  $\text{StreamAlg}$ .
- $\text{StreamAlg.Update}(x_t)$  : On input an additive update  $x_t := (i_t, \Delta_t)$ , the update operation modifies  $\mathcal{DS}$ .

- $\text{StreamAlg.Report}() \rightarrow r$ : The report operation uses  $\mathcal{DS}$  to compute a query response  $r \in X$ .

A response  $r$  is said to be correct if  $r = f_n(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{Z}^n$  is the accumulation of all past updates such that  $\mathbf{x}_i = \sum_t \mathbf{1}(i_t = i) \cdot \Delta_t$  for  $i \in [n]$ .

In applications,  $\text{StreamAlg.Setup}$  is usually called once before the beginning of the stream, and  $\text{StreamAlg.Update}$  is called once on each update  $x_t$  during the stream. Depending on the use case,  $\text{StreamAlg.Report}$  may be called once or multiple times to get query responses.

In this work, we only consider exact computations of query functions ( $r = f_n(\mathbf{x})$ ), though the above definition can also be generalized to approximation problems ( $r \approx f_n(\mathbf{x})$ ).

*White-box Adversarial Stream.* In the conventional *oblivious* stream model, the sequence of updates maybe chosen adversarially, but is assumed to be chosen independently of the randomness of a streaming algorithm  $\text{StreamAlg}$ . In contrast, in the white-box adversarial streaming model, a sequence of stream updates  $(x_1, \dots, x_T)$  is adaptively generated by an adversary who sees the *full internal state* of  $\text{StreamAlg}$  at all times, including  $\mathcal{DS}$  and any other parameters and randomness used by  $\text{StreamAlg}$ . We say that a streaming algorithm is *white-box adversarially robust* if it guarantees to output correct responses even against such powerful adversaries. This is described as the following game:

**Definition 2.12 (White-Box Adversarially Robustness).** We say a streaming algorithm  $\text{StreamAlg}$  is white-box adversarially robust (WAR) if for all dimension parameters  $n \in \mathbb{N}$ , the following holds: For any PPT adversary  $\mathcal{A}$ , the following experiment  $\text{Expt}_{\mathcal{A}, \text{StreamAlg}}(n)$  outputs 1 with probability at most  $\text{negl}(n)$ :

$\text{Expt}_{\mathcal{A}, \text{StreamAlg}}(n)$  :

1. The challenger and  $\mathcal{A}$  agree on a query function ensemble  $\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow X\}_{n \in \mathbb{N}}$  for some domain  $X$ . On input a dimension parameter  $n$ , let  $f_n \in \mathcal{F}$  be the query function of the experiment.
2. The challenger runs  $\text{StreamAlg.Setup}(n)$ . It then provides all internal states of  $\text{StreamAlg}$  to  $\mathcal{A}$ , including  $\mathcal{DS}$  and any other parameters and randomness previously used by  $\text{StreamAlg}$ .
3. For some  $T \in \text{poly}(n)$ , at each time  $t \in [T]$ :
  - $\mathcal{A}$  outputs an adaptive update  $x_t := (i_t, \Delta_t)$  for  $t \in [T]$ . Also,  $\mathcal{A}$  may issue a query.
  - The challenger runs  $\text{StreamAlg.Update}(x_t)$ , and outputs a response  $r_t \leftarrow \text{StreamAlg.Report}()$  if it is queried.

The challenger then provides all internal states of  $\text{StreamAlg}$  to  $\mathcal{A}$ , including  $\mathcal{DS}$  and any other parameters and randomness used by  $\text{StreamAlg}$ .
4. The experiment outputs 1 if and only if at some time  $t \in [T]$  in the stream,  $\text{StreamAlg}$  output an incorrect response  $r_t \neq f_n(\mathbf{x})$ , where  $\mathbf{x}$  is the underlying data vector with  $\mathbf{x}_i = \sum_{j \in [t]} \mathbf{1}(i_j = i) \cdot \Delta_j$  for  $i \in [n]$ .



### 3 Streaming $K$ -Sparse Recovery

*Sparse Recovery Problem* We study the sparse recovery problem for vectors in white-box adversarial streaming settings. There exist deterministic (thus WAR) streaming algorithms that recover  $k$ -sparse vectors *assuming that the input vector is  $k$ -sparse*. However, their outputs can be arbitrary when the input violates such an assumption. In contrast, our algorithm has the crucial property that it can detect if the input violates the sparsity constraint. We describe the problem as follows:

**Definition 3.1 ( $k$ -Sparse Recovery).** *Given a sparsity parameter  $k$ , the query function ensemble of the  $k$ -Sparse Recovery problem is:*

$$\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow (\mathbb{Z}^n \cup \{\perp\})\}_{n \in \mathbb{N}}, \text{ where } f_n(\mathbf{x}) = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_0 \leq k \\ \perp & \text{otherwise} \end{cases}$$

Here and throughout, for the sparsity to make sense we always assume  $k \leq n$ . And  $k$  can be a function in  $n$ .

Before tackling this problem, we first consider a relaxation for it, which only guarantees correct recovery when the input vector  $\mathbf{x}$  is  $k$ -sparse. In other words, we have no constraint on the function output when  $\|\mathbf{x}\|_0 > k$ .

**Definition 3.2 (Relaxed  $k$ -Sparse Recovery).** *Given a sparsity parameter  $k$ ,  $\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow \mathbb{Z}^n\}_{n \in \mathbb{N}}$  is a query function ensemble of the Relaxed  $k$ -Sparse Recovery problem if for all  $f_n \in \mathcal{F}$ , the following holds: For all  $\mathbf{x} \in \mathbb{Z}^n$  satisfying  $\|\mathbf{x}\|_0 \leq k$ ,  $f_n(\mathbf{x}) = \mathbf{x}$ .*

As mentioned earlier, there exists space- and time-efficient, deterministic streaming algorithms that solve the relaxed  $k$ -sparse recovery problem. We restate the existing result as follows, which we will use as a black box subroutine in our algorithm.

**Theorem 3.3 ([44,43,34]).** *There exists a streaming algorithm  $\text{StreamAlg}_0$  that solves the Relaxed  $k$ -Sparse Recovery problem (Definition 3.2), with  $\text{StreamAlg}_0.\text{Setup}$ ,  $\text{Update}$ , and  $\text{Report}$  all deterministic and run in  $\text{poly}(n)$  time.*

We will instantiate  $\text{StreamAlg}_0$  as concrete existing schemes and specify their space and time complexity when discussing the efficiency of our construction in Section 4. But for the purpose of proving that our construction is WAR, we only need the polynomial runtime and the deterministic property.

With the result in Theorem 4.1 and a PFHE, we are now ready to construct a WAR streaming algorithm for the (strong)  $k$ -sparse recovery problem.

**Construction 3.1.** *Let  $\text{StreamAlg}_0$  be an algorithm for the Relaxed  $k$ -Sparse Recovery problem from Theorem 4.1. And let PFHE be a pseudorandom FHE satisfying Definition 2.10.*

*Our construction is as follows:*

- **StreamAlg.Setup( $n$ ):** Run **StreamAlg<sub>0</sub>.Setup( $n$ )**.  
Define  $\ell = \lceil \log n \rceil$ . For each  $i \in [n]$ , define a circuit  $C_i : \{0,1\}^\ell \rightarrow \{0,1\}$ , such that:

$$C_i(\mu_1, \mu_2, \dots, \mu_\ell) = \begin{cases} 1 & \text{if } (\mu_1, \mu_2, \dots, \mu_\ell) \text{ is the bit representation of } i \\ 0 & \text{otherwise.} \end{cases}$$

Let  $\lambda = \mathcal{S}(n)$  be the security parameter of PFHE, for some function  $\mathcal{S}$  in  $n$  that we will specify later in Section 4. Choose a  $q \in \text{poly}(n)$  with  $q \gg n \cdot N$  to be the modulus of the space of PFHE ciphertexts, where  $N \in \text{poly}(n)$  is the bound of the stream. All computations will be modulo  $q$ .

Take  $\mathcal{D}_{\widetilde{\text{pk}}(\lambda)}$  to be the distribution of pseudo-public key in Definition 2.10.

Sample  $\widetilde{\text{pk}} \leftarrow \mathcal{D}_{\widetilde{\text{pk}}(\lambda)}$ .

Similarly, take  $\mathcal{D}_{\widetilde{\text{ct}}(\lambda)}$  to be the distribution of pseudo-ciphertext in Definition 2.10. Sample  $\widetilde{\text{ct}}_1, \widetilde{\text{ct}}_2, \dots, \widetilde{\text{ct}}_\ell \leftarrow \mathcal{D}_{\widetilde{\text{ct}}(\lambda)}$ . Store  $(\widetilde{\text{pk}}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_\ell)$  and  $\text{sketch} := 0$  in the internal data structure  $\mathcal{DS}$  of **StreamAlg**.

- **StreamAlg.Update( $x_t$ ):** Run **StreamAlg<sub>0</sub>.Update( $x_t$ )**.  
For  $x_t := (i_t, \Delta_t)$ , evaluate  $\widehat{\text{ct}}_{i_t} \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_{i_t}, \widetilde{\text{ct}}_1, \widetilde{\text{ct}}_2, \dots, \widetilde{\text{ct}}_\ell)$ . Then update sketch as

$$\text{sketch} \leftarrow \text{sketch} + \Delta_t \cdot \widehat{\text{ct}}_{i_t}.$$

- **StreamAlg.Report():** Let  $\mathbf{x}' \leftarrow \text{StreamAlg}_0.\text{Report}()$ . If  $\|\mathbf{x}'\|_0 > k$ , output  $\perp$ .  
Otherwise, evaluate  $\widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_i, \widetilde{\text{ct}}_1, \widetilde{\text{ct}}_2, \dots, \widetilde{\text{ct}}_\ell)$  for  $i \in [n]$ . Then compute

$$\text{hash} \leftarrow \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot \mathbf{x}'_i.$$

Output  $\mathbf{x}'$  if  $\text{hash} = \text{sketch}$ , and output  $\perp$  otherwise.

**Robustness.** Recall that in the FHE Preliminary section, we define the  $(\mathcal{T}, \epsilon)$ -security of a FHE scheme in Remark 2.9, where  $\mathcal{T}$  and  $\epsilon$  are functions in  $\lambda$ . Also, in Construction 3.1, we set  $\lambda = \mathcal{S}(n)$  for some function  $\mathcal{S}$ .

In this section, we will show that given appropriate relationships between these functions, the **StreamAlg** scheme in Construction 3.1 is WAR. We state our main theorem as follows:

**Theorem 3.4 (WAR of Construction 3.1).** *Suppose PFHE is a  $(\mathcal{T}, \epsilon)$ -pseudorandom FHE scheme from Definition 2.10 with  $(\mathcal{T}, \epsilon)$ -security, and **StreamAlg<sub>0</sub>** is a streaming algorithm for Relaxed  $k$ -Sparse Recovery from Theorem 4.1. If for all  $n \in \mathbb{N}$ , functions  $\mathcal{T}, \epsilon$ , and  $\mathcal{S}$  satisfy:*

- $\mathcal{T}(\mathcal{S}(n)) = \Omega(\text{poly } n)$
- $\epsilon(\mathcal{S}(n)) \leq \text{negl}(n)$  for some negligible functions  $\text{negl}$ .

*then the **StreamAlg** scheme in Construction 3.1 is a WAR streaming algorithm that solves the  $k$ -Sparse Recovery problem (Definition 3.1).*

Roughly speaking, we need  $\text{poly}(n) \subseteq \mathcal{T}(\mathcal{S}(n))$  to convert between the runtime of a PFHE adversary and a **StreamAlg** adversary. And we need  $\epsilon(\mathcal{S}(n)) \leq \text{negl}(n)$  to convert between the advantage of a PFHE adversary and a **StreamAlg** adversary.

*Proof.* We first show that at any time during the stream, if the underlying vector  $\mathbf{x}$  is  $k$ -sparse, then the query response of **StreamAlg** is correct.

**Lemma 3.5.** *Suppose PFHE satisfies Definition 2.10 and **StreamAlg**<sub>0</sub> solves the Relaxed  $k$ -Sparse Recovery problem, the following holds: At any time  $t \in [T]$ , if the current underlying vector  $\mathbf{x}$  satisfies  $\|\mathbf{x}\|_0 \leq k$ , then **StreamAlg**.Report correctly outputs  $\mathbf{x}$ .*

*Proof.* By the correctness of the **StreamAlg**<sub>0</sub> scheme, when  $\|\mathbf{x}\|_0 \leq k$ , **StreamAlg**<sub>0</sub>.Report outputs  $\mathbf{x}' = \mathbf{x}$ . Since Definition 2.10 restricts PFHE.Eval to be deterministic, we essentially have  $\text{hash} = \sum_{i \in [n]} \hat{c}t_i \cdot \mathbf{x}'_i$  and  $\text{sketch} = \sum_{i \in [n]} \hat{c}t_i \cdot \mathbf{x}_i$  for the same set of ciphertexts  $\hat{c}t_i$ . Thus  $\text{hash} = \text{sketch}$  follows from  $\mathbf{x}' = \mathbf{x}$ . In this case, **StreamAlg**.Report correctly outputs  $\mathbf{x}$  as specified by the query function  $f_n$  of the  $k$ -Sparsed Recovery problem.  $\square$

It remains to show the correctness of responses in the case when the underlying data vector is denser than  $k$ . Our proof proceeds via a sequence of hybrid experiments between an adversary  $\mathcal{A}$  and a challenger.

- **Hybrid**<sub>0</sub>: This is a real WAR experiment  $\text{Expt}_{\mathcal{A}, \text{StreamAlg}}(n)$  from Definition 2.12, where the challenger runs the **StreamAlg** scheme from Construction 3.1. On input a dimension parameter  $n$ , the challenger samples  $\tilde{\mathbf{pk}} \leftarrow \mathcal{D}_{\tilde{\mathbf{pk}}(\lambda)}$  and  $\tilde{\mathbf{ct}}_1, \tilde{\mathbf{ct}}_2, \dots, \tilde{\mathbf{ct}}_{\lceil \log n \rceil} \leftarrow \mathcal{D}_{\tilde{\mathbf{ct}}(\lambda)}$  from the distributions of pseudo-public key and pseudo-ciphertext (Definition 2.10) and provides them to  $\mathcal{A}$ .  $\mathcal{A}$  outputs  $T \in \text{poly}(n)$ -many adaptive updates  $x_t := (i_t, \Delta_t)$  upon seeing past internal states of **StreamAlg**. At any time  $t \in [T]$ , these updates implicitly define an underlying vector  $\mathbf{x}$ , whose  $i$ -th coordinate  $\mathbf{x}_i = \sum_t \mathbf{1}(i_t = i) \cdot \Delta_t$  for  $i \in [n]$ .
- **Hybrid**<sub>1</sub>: This is the same experiment as **Hybrid**<sub>0</sub>, except that  $\mathcal{A}$  now runs some additional extraction on itself:
  - At each time  $t \in [T]$ ,  $\mathcal{A}$  computes the underlying vector  $\mathbf{x}$  defined by all of its past updates.
  - $\mathcal{A}$  runs an additional instance of the **StreamAlg**<sub>0</sub> scheme from Definition 3.2 (which solves the Relaxed  $k$ -Sparse Recovery problem):
    1. Before  $\mathcal{A}$  generates its first update, it initializes its **StreamAlg**<sub>0</sub> instance to be in the same state as the **StreamAlg**<sub>0</sub> scheme run by the challenger. (i.e.,  $\mathcal{A}$  runs **StreamAlg**<sub>0</sub>.Setup( $n$ ) using the same set of randomness used by the challenger.)
    2. At each time  $t \in [T]$ , after  $\mathcal{A}$  generates an update  $x_t$ , it calls **StreamAlg**<sub>0</sub>.Update( $x_t$ ) and then **StreamAlg**<sub>0</sub>.Report() to get a query response  $\mathbf{x}'$ .

We note that if  $\mathcal{A}$  runs in  $\text{poly}(n)$  time in **Hybrid**<sub>0</sub>, then it still runs in  $\text{poly}(n)$  time in **Hybrid**<sub>1</sub> since computing  $\mathbf{x}$  is trivial and **StreamAlg**<sub>0</sub> runs in poly-time as specified in Theorem 4.1.

For an adversary  $\mathcal{A}$ , we write  $\text{Hyb}_i(\mathcal{A})$  to denote the output of **Hybrid** <sub>$i$</sub> .

**Lemma 3.6.** *For all adversaries  $\mathcal{A}$ ,  $|\Pr[\text{Hyb}_0(\mathcal{A}) = 1] - \Pr[\text{Hyb}_1(\mathcal{A}) = 1]| = 0$ .*

*Proof.* This just follows from the fact that the additional extractions run by  $\mathcal{A}$  do not affect the output of the experiments.  $\square$

– **Hybrid**<sub>2</sub>: Same as **Hybrid**<sub>1</sub>, except that at the beginning of the experiment, the challenger samples a random coordinate  $m \xleftarrow{\$} [n]$  which will be hidden from  $\mathcal{A}$ . The experiment outputs 1 if at some time  $t \in [T]$  both of the following holds:

1. **StreamAlg** outputs an incorrect response  $r_t \neq f_n(\mathbf{x})$  (i.e. same condition for  $\text{Hyb}_1(\mathcal{A}) = 1$ ),
2. The  $m$ -th coordinates of  $\mathbf{x}$  and  $\mathbf{x}'$  extracted by  $\mathcal{A}$  satisfy  $x_m \neq x'_m$ ,

We show that with the additional success condition, the success probability of **Hybrid**<sub>2</sub> drops by at most a  $1/n$  multiplicative factor compared to **Hybrid**<sub>1</sub>.

**Lemma 3.7.** *For all adversaries  $\mathcal{A}$ ,  $\Pr[\text{Hyb}_2(\mathcal{A}) = 1] \geq \Pr[\text{Hyb}_1(\mathcal{A}) = 1]/n$ .*

Before proving Lemma 3.7, we make the following observation:

*Claim.* At any time  $t \in [T]$  during the **Hybrid**<sub>2</sub> experiment, if **StreamAlg** outputs an incorrect response  $r_t \neq f_n(\mathbf{x})$ , and  $\mathbf{x}$  and  $\mathbf{x}'$  are the vectors extracted by  $\mathcal{A}$  at time  $t$ , then  $\mathbf{x}' \neq \mathbf{x}$ .

*Proof.* By Lemma 3.5,  $r_t \neq f_n(\mathbf{x})$  can occur only when  $\|\mathbf{x}\|_0 > k$  and  $r_t \neq \perp$ . In this case, by the construction of **StreamAlg.Report**,  $r_t$  should be a  $k$ -sparse vector thus different from  $\mathbf{x}$ . Also, using deterministic **StreamAlg**<sub>0</sub>.Update and **StreamAlg**<sub>0</sub>.Report,  $\mathcal{A}$  exactly extracts  $\mathbf{x}' = r_t$ , so  $\mathbf{x}' \neq \mathbf{x}$ .  $\square$

Lemma 3.7 then follows from Claim 3:

*Proof (of Lemma 3.7).* Claim 3 states that  $\mathbf{x}' \neq \mathbf{x}$ , and in particular,  $\mathbf{x}'$  and  $\mathbf{x}$  disagree at at least one coordinate. Recall that the random coordinate  $m$  is hidden from  $\mathcal{A}$  at all times, so with probability at least  $1/n$ ,  $m$  collides with a coordinate where  $\mathbf{x}'$  and  $\mathbf{x}$  disagrees. In this case, both conditions in **Hybrid**<sub>2</sub> are satisfied, so  $\Pr[\text{Hyb}_2(\mathcal{A}) = 1 \mid \text{Hyb}_1(\mathcal{A}) = 1] \geq 1/n$ .  $\square$

– **Hybrid**<sub>3</sub>: Same as **Hybrid**<sub>2</sub>, except that the challenger now generates an actual public key  $\text{pk} : (\text{pk}, \text{sk}) \leftarrow \text{PFHE.Setup}(1^\lambda)$  to replace  $\widetilde{\text{pk}}$ .

– **Hybrid**<sub>4</sub>: Same as **Hybrid**<sub>3</sub>, except that the challenger now uses the random coordinate  $m \xleftarrow{\$} [n]$  (whose bits are denoted as  $m_1 m_2 \dots m_{\lceil \log n \rceil}$ ) that was sampled privately as the message in the ciphertexts: It generates  $\lceil \log n \rceil$  actual ciphertexts  $\text{ct}_i \leftarrow \text{PFHE.Enc}(\text{pk}, m_i)$  for  $i \in [\lceil \log n \rceil]$  to replace  $(\text{ct}_1, \text{ct}_2, \dots, \text{ct}_{\lceil \log n \rceil})$ . i.e.,  $(\text{ct}_1, \text{ct}_2, \dots, \text{ct}_{\lceil \log n \rceil})$  encrypts the bit representation of  $m$ .

**Lemma 3.8.** *Given that PFHE satisfies Definition 2.10, for any  $PPT_n$  adversaries  $\mathcal{A}$ ,  $|\Pr[\text{Hyb}_2(\mathcal{A}) = 1] - \Pr[\text{Hyb}_3(\mathcal{A}) = 1]| \leq \text{negl}(n)$  for some negligible function  $\text{negl}$ .*

*Proof.* The only difference in  $\mathcal{A}$ 's views in **Hybrid**<sub>2</sub> and **Hybrid**<sub>3</sub> is the way the challenger generates  $\widetilde{\text{pk}}$  or  $\text{pk}$ . In **Hybrid**<sub>2</sub>, the challenger samples  $\widetilde{\text{pk}} \leftarrow \mathcal{D}_{\widetilde{\text{pk}}(\lambda)}$ ; and in **Hybrid**<sub>3</sub>, the challenger runs  $\text{pk} : (\text{pk}, \text{sk}) \leftarrow \text{PFHE.Setup}(1^\lambda)$ . Thus  $\mathcal{A}$  distinguishing these two hybrids will distinguish the two distributions  $\{\text{pk} : (\text{pk}, \text{sk}) \leftarrow \text{FHE.Setup}(1^\lambda)\}$  and  $\mathcal{D}_{\widetilde{\text{pk}}(\lambda)}$  with the same advantage. And by the  $(\mathcal{T}, \epsilon)$ -indistinguishability in Definition 2.10, for any  $\text{poly}(n) \leq \mathcal{T}(\lambda)$ -sized adversaries, the latter two distributions can be distinguished with probability at most  $\epsilon(\lambda) \leq \text{negl}(n)$ .  $\square$

**Lemma 3.9.** *Given that PFHE satisfies Definition 2.10, for any  $PPT_n$  adversaries  $\mathcal{A}$ ,  $|\Pr[\text{Hyb}_3(\mathcal{A}) = 1] - \Pr[\text{Hyb}_4(\mathcal{A}) = 1]| \leq \text{negl}(n)$  for some negligible function  $\text{negl}$ .*

*Proof.* Similar to Lemma 3.8, the only difference in  $\mathcal{A}$ 's views in **Hybrid**<sub>3</sub> and **Hybrid**<sub>4</sub> is the way the challenger generates  $(\widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  or  $(\text{ct}_1, \dots, \text{ct}_{\lceil \log n \rceil})$ . Recall that we have:  $\{\text{ct}\} \approx_{(\mathcal{T}, \epsilon)} \mathcal{D}_{\widetilde{\text{ct}}(\lambda)}$  from Definition 2.10. Therefore for any  $PPT_n$  adversary  $\mathcal{A}$ , the advantage of distinguishing  $(\widetilde{\text{ct}}_1, \widetilde{\text{ct}}_2, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  in **Hybrid**<sub>3</sub> from a ciphertext sequence  $(\text{ct}_1, \text{ct}_2, \dots, \text{ct}_{\lceil \log n \rceil})$  where each  $\text{ct}_i \leftarrow \text{PFHE.Enc}(\text{pk}, m_i)$  is at most  $\lceil \log n \rceil \cdot \epsilon(\lambda) = \text{negl}(n)$ . Thus  $\mathcal{A}$  cannot distinguish between **Hybrid**<sub>3</sub> and **Hybrid**<sub>4</sub> with non-negligible advantage, as otherwise either the pseudorandom-ct property in Definition 2.10 would be broken.  $\square$

Now that we have related the success probability between the above consecutive hybrid experiments, we will show next that  $\Pr[\text{Hyb}_4(\mathcal{A}) = 1] = 0$ .

**Lemma 3.10.** *Given that PFHE satisfies Definition 2.10, for any  $PPT_n$  adversaries  $\mathcal{A}$ ,  $\Pr[\text{Hyb}_4(\mathcal{A}) = 1] = 0$ .*

*Proof.* As in the proof for Claim 3, if  $\text{Hyb}_4(\mathcal{A}) = 1$ , at some time  $t$  we must have  $r_t = \mathbf{x}' \neq \mathbf{x}$ . Additionally, we have:

- $\mathbf{x}'_m \neq \mathbf{x}_m$ , which is among the condition of  $\text{Hyb}_4(\mathcal{A}) = 1$ , and
- $\text{hash} = \text{sketch}$ , which is verified in  $\text{StreamAlg.Report}$ .

Recall the  $\text{StreamAlg}$  scheme computes  $\text{hash} := \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot \mathbf{x}'_i$  and  $\text{sketch} := \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot \mathbf{x}_i$ , where  $\widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\text{pk}, C_i, \text{ct}_1, \dots, \text{ct}_{\lceil \log n \rceil})$ , for  $C_i(\mu) = 1$  iff  $\mu = i$ . Therefore, given that  $(\text{ct}_1, \dots, \text{ct}_{\lceil \log n \rceil})$  encrypts  $m$  in **Hybrid**<sub>4</sub>,  $\widehat{\text{ct}}_m$  should encode 1, and  $\widehat{\text{ct}}_i$  encodes 0 for all  $i \neq m$ . Thus by the linear homomorphism property of PFHE, there exists a deterministic decoding algorithm such that  $\text{hash}$  decodes to  $\mathbf{x}'_m$  and  $\text{sketch}$  decodes to  $\mathbf{x}_m$ . However, since  $\mathbf{x}'_m \neq \mathbf{x}_m$ , this essentially implies  $\text{hash} \neq \text{sketch}$ , contradicting the condition in  $\text{StreamAlg.Report}$ .  $\square$

Finally, summarizing Lemma 3.6 to Lemma 3.10:

1.  $|\Pr[\text{Hyb}_0(\mathcal{A}) = 1] - \Pr[\text{Hyb}_1(\mathcal{A}) = 1]| = 0$
2.  $\Pr[\text{Hyb}_2(\mathcal{A}) = 1] \geq \Pr[\text{Hyb}_1(\mathcal{A}) = 1]/n$
3.  $|\Pr[\text{Hyb}_2(\mathcal{A}) = 1] - \Pr[\text{Hyb}_3(\mathcal{A}) = 1]| = \text{negl}(n)$
4.  $|\Pr[\text{Hyb}_3(\mathcal{A}) = 1] - \Pr[\text{Hyb}_4(\mathcal{A}) = 1]| = \text{negl}(n)$
5.  $\Pr[\text{Hyb}_4(\mathcal{A}) = 1] = 0$

We have that for any adversaries  $\mathcal{A}$ ,  $\Pr[\text{Hyb}_0(\mathcal{A}) = 1] \leq \text{negl}(n) \cdot n = \text{negl}(n)$ . i.e., The real WAR game between  $\mathcal{A}$  and a challenger running the **StreamAlg** scheme from Construction 3.1 outputs 1 with probability at most  $\text{negl}(n)$ .  $\square$

## 4 Efficiency

In this section, we analyze the bit complexity of the **StreamAlg** scheme in Construction 3.1, in addition to the runtime of **StreamAlg.Update** and **StreamAlg.Report** (which we refer to as *update time* and *report time*, respectively). We start by instantiating **StreamAlg<sub>0</sub>** and PFHE used in our construction.

*Notation.* We use  $\mathbf{b}_i$  to denote the standard basis vector, whose  $i$ -th entry equals 1 and 0 anywhere else. Our logarithms are in base 2.

### 4.1 Deterministic Relaxed Sparse Recovery

Recall that in Theorem 4.1, we state that there exists space- and time-efficient streaming algorithms for  $k$ -sparse recovery without detection (Definition 3.2). In general, given a sparse vector  $\mathbf{x} \in \mathbb{Z}^n$ , these algorithms perform some *linear measurement*  $\Phi : \mathbb{Z}^n \rightarrow \mathbb{Z}^r$  on  $\mathbf{x}$  with  $r \ll n$  to get a compressive representation  $\Phi(\mathbf{x})$ , which is often referred to as a *sketch* of  $\mathbf{x}$ . This sketch is later used to reconstruct  $\mathbf{x}$ . The linearity property is useful for incrementally maintaining the sketch during the stream. For example, updating the sketch after incrementing the  $i$ -th coordinate  $\mathbf{x}_i$  can be done as  $\Phi(\mathbf{x} + \mathbf{b}_i) = \Phi(\mathbf{x}) + \Phi(\mathbf{b}_i)$ . When describing  $\Phi$  as a matrix, we implicitly define the mapping of a vector  $\mathbf{x}$  as the matrix-vector product  $\Phi\mathbf{x}$ .

Many existing algorithms implement the above specification using a random measurement matrix  $\Phi$ . However, in our attempt to remove the random oracle assumption within sub-linear total space, we cannot afford to store an  $n$ -column random matrix during the stream. To this end, we use algorithms with deterministic construction, such that the measurements themselves incur no space overhead. We summarize two existing results and discuss their applicability for our purpose.

**Theorem 4.1 ([34]).** *There exists a streaming algorithm **StreamAlg<sub>0</sub>** for  $k$ -Sparse Recovery without Detection, such that given an input parameter  $n$ , for an integer stream with entries bounded by  $\text{poly}(n)$ :*

- $\text{StreamAlg}_0.\text{Setup}, \text{Update}, \text{and Report}$  are all deterministic.
- $\text{StreamAlg}_0$  takes  $\tilde{O}(k)$  bits of space.
- $\text{StreamAlg}_0.\text{Update}$  runs in  $\tilde{O}(1)$ .
- $\text{StreamAlg}_0.\text{Report}$  runs in  $\tilde{O}(k^{1+c})$  for an arbitrarily small constant  $c > 0$ .

The algorithm in Theorem 4.1 uses a two-layer hashing as its measurement, which is capable of recovering all length- $n$ ,  $k$ -sparse vectors  $\mathbf{x}$ . It achieves nearly optimal space and time. Though such measurement is “constructed” non-explicitly using the probabilistic method. For explicit construction, we consider the following result:

**Theorem 4.2 ([44]).** *There exists a streaming algorithm  $\text{StreamAlg}_0$  for  $k$ -Sparse Recovery without Detection, such that given an input parameter  $n$ , for an integer stream with entries bounded by  $\text{poly}(n)$ :*

- $\text{StreamAlg}_0.\text{Setup}, \text{Update}, \text{and Report}$  are all deterministic.
- $\text{StreamAlg}_0$  takes  $\tilde{O}(k)$  bits of space.
- $\text{StreamAlg}_0.\text{Update}$  runs in  $\tilde{O}(1)$ .
- $\text{StreamAlg}_0.\text{Report}$  runs in  $\tilde{O}(k^2)$ .

The measurement  $\Phi$  in Theorem 4.2 is a  $2k \times n$  Vandermonde matrix explicitly defined<sup>1</sup> as  $\Phi_{i,j} = j^{i-1}$  over  $\mathbb{Z}_p$ , where  $p > 2M$  is a large enough prime for  $M \in \text{poly}(n)$  being the bound on the magnitude of the stream entries. The reconstruction algorithm uses the same idea as the algebraic algorithm for decoding Reed-Solomon codes. It uses the sketch to construct an error-locator polynomial; the roots of this polynomial identify the signals appearing in the sparse superposition.

We note that the Vandermonde matrix is dense but structural. Thus to achieve the claimed  $\tilde{O}(1)$  update time, we buffer every  $2k$  stream updates and then efficiently process them in one batch. For preliminaries, we need the following result for fast multi-point evaluation of polynomials and the Transposition Principle:

**Theorem 4.3 ([32]).** *Let  $R$  be a ring, and let  $q \in R[x]$  be a degree- $d$  polynomial. Then given at most  $d$  distinct points  $x_1, \dots, x_d \in R$ , all values  $q(x_1), \dots, q(x_d)$  can be computed using  $\mathcal{O}(d \log^2 d \log \log d)$  total operations over  $R$ .*

**Theorem 4.4 (Transposition Principle, [63]).** *Let  $K$  be a field and  $\mathbf{A}$  be an  $r \times w$  matrix over  $K$ . Suppose there exists an arithmetic circuit of size  $s$  that can compute  $\mathbf{A}\mathbf{y}$  for arbitrary length- $w$  vector  $\mathbf{y}$  over  $K$ . Then there exists an  $\mathcal{O}(s)$ -time algorithm that transforms this circuit to compute  $\mathbf{A}^\top \mathbf{x}$  for arbitrary length- $r$  vector  $\mathbf{x}$  over  $K$ , with the size of the new circuit in  $\mathcal{O}(s)$ .*

Theorem 4.3 and Theorem 4.4 together show the following update time for Vandermonde measurements:

<sup>1</sup>For our discussion, we choose a Vandermonde matrix that is easy to describe. Though many other Vandermonde matrices also satisfy the specifications.

**Proposition 4.5.** *Suppose  $\Phi$  is a  $2k \times n$  Vandermonde matrix as described above. Let  $x_1, x_2, \dots, x_{2k}$  be a sequence of updates with  $x_t := (i_t, \Delta_t)$  for each  $t \in [2k]$ . Define  $\mathbf{v} = \sum_{t \in [2k]} \mathbf{b}_{i_t} \cdot \Delta_t$ , then  $\Phi(\mathbf{v})$  can be computed using  $\tilde{O}(k)$  time.*

*Proof.* Without loss of generality, we assume all updated coordinates  $i_1, i_2, \dots, i_{2k}$  are distinct, as one can always compress multiple updates to the same coordinate as a single update. Then  $i_1, i_2, \dots, i_{2k}$  index  $2k$  columns of  $\Phi$  and  $2k$  coordinates of  $\mathbf{v}$  (which are the only coordinates that can possibly be non-zero). Let  $\mathbf{A}$  be the indexed  $2k \times 2k$  column sub-matrix and  $\mathbf{v}'$  be the  $2k$ -length vector containing only  $i_1, i_2, \dots, i_{2k}$ -th entries of  $\mathbf{v}$  in their original order. We have  $\Phi \mathbf{v} = \mathbf{A} \mathbf{v}'$ .

Observe that the column submatrix  $\mathbf{A}$  is still a Vandermonde matrix, and each row  $i$  of  $\mathbf{A}^\top$  takes the form  $[1, \alpha_i, \alpha_i^2, \dots, \alpha_i^{2k-1}]$  for some  $\alpha_i \in [n]$ . For an arbitrary vector  $\mathbf{u} \in \mathbb{Z}_p^{2k}$ , we interpret it as the coefficient vector of a degree- $2k$  polynomial in  $\mathbb{Z}_p^{2k}$ , then  $\mathbf{A}^\top \mathbf{u}$  can be computed as evaluating the polynomial at  $2k$  points  $\alpha_1, \alpha_2, \dots, \alpha_{2k}$ , which takes  $\tilde{O}(k)$  time using the fast multi-point evaluation in Theorem 4.1. By the Transposition Principle, this implies asymptotically the same runtime for evaluating  $\mathbf{A} \mathbf{v}'$  thus  $\Phi \mathbf{v}$ .  $\square$

Therefore, once we buffer the updates into  $2k$ -sized batches and amortize the cost of processing each batch to the future  $2k$  updates, we get  $\tilde{O}(1)$  amortized runtime for `StreamAlg0.Update`.

## 4.2 Pseudorandom FHE

We focus on GSW [33], a standard LWE based fully homomorphic encryption scheme, for our instantiation. In this section, we first show that GSW has the additional properties that we require for a PFHE (Definition 2.10). Then we specify our security parameter  $\lambda$ , assumption, and functions  $(\mathcal{T}, \epsilon)$ , and state the space and time complexity of PFHE given these parameters.

*GSW Satisfies PFHE Properties.* To start with, we briefly review the construction of GSW. We omit Dec since it is not involved in our scheme. Instead, we will describe a decoding algorithm `LinearDec` for linear combinations of ciphertexts (which largely resembles the original Dec of GSW) when showing the  $\mathbb{Z}$  Linear Homomorphism property of GSW.

- `Setup( $1^\lambda$ )`: Choose a modulus  $q$ , a lattice dimension parameter  $g$ , and error distribution  $\chi$  parametrized by  $\lambda$ . Also, choose some  $h = \Theta(g \log q)$ . Let  $params = (g, h, q, \chi)$  be parameters for  $\text{LWE}_{g,h,q,\chi}$ .
- `SKGen( $params$ )`: Sample  $\bar{\mathbf{s}} \leftarrow \mathbb{Z}_q^g$ . Output  $\mathbf{sk} = \mathbf{s}^\top = [-\bar{\mathbf{s}}^\top \mid 1]$ .
- `PKGen( $params, \mathbf{sk}$ )`: Sample a matrix  $\bar{\mathbf{A}} \leftarrow \mathbb{Z}_q^{g \times h}$  uniformly and an error vector  $\mathbf{e} \leftarrow \chi^h$ . Let  $\boldsymbol{\alpha} = \bar{\mathbf{s}}^\top \cdot \bar{\mathbf{A}} + \mathbf{e}^\top$ . Set  $\mathbf{pk} = \mathbf{A}$  to be the  $(g+1)$ -row matrix consisting of  $g$  rows of  $\bar{\mathbf{A}}$  followed by  $\boldsymbol{\alpha}$ .
- `Enc( $params, \mathbf{pk}, \mu$ )`: Sample a uniform matrix  $\mathbf{R} \leftarrow \{0, 1\}^{h \times h}$ . For  $\mathbf{pk} = \mathbf{A}$ , output  $\mathbf{A} \cdot \mathbf{R} + \mu \cdot \mathbf{G}$  where  $\mathbf{G}$  is a gadget matrix.



Similar to the previous Remark 2.9, we assume that our set of LWE parameters  $params$  is  $(\mathcal{T}, \epsilon)$ -secure, such that for any  $\mathcal{T}(\lambda)$  time adversary, the advantage of distinguishing an LWE sample from a uniformly random sample is at most  $\epsilon(\lambda)$ . We now check that GSW satisfies the additional properties required by PFHE:

- Eval is deterministic by the construction of GSW.
- (Pseudorandom pk, ct.) This is just the built-in security of GSW shown via the LWE hardness and the Left-over Hash lemma. Assuming the  $(\mathcal{T}, \epsilon)$ -hardness of  $\text{LWE}_{g,h,q,\chi}$ , we have:

$$\left\{ \begin{array}{l} \mathbf{A} \leftarrow \text{PKGen} \\ (\mathbf{A}, \mathbf{A} \cdot \mathbf{R} + \mu \cdot \mathbf{G}) : \mathbf{R} \leftarrow \{0, 1\}^{h \times h} \\ \mu \in \{0, 1\} \end{array} \right\} \approx_{(\mathcal{T}, \epsilon)} \left\{ (\mathbf{U}_1, \mathbf{U}_2) : \begin{array}{l} \mathbf{U}_1 \leftarrow \mathbb{Z}_q^{(g+1) \times h} \\ \mathbf{U}_2 \leftarrow \mathbb{Z}_q^{(g+1) \times h} \end{array} \right\}.$$

- ( $\mathbb{Z}$  Linear Homomorphism.) Consider the following arbitrary linear combination, where each ciphertext  $\text{ct}_i = \mathbf{A} \cdot \mathbf{R}_i + \mu_i \cdot \mathbf{G}$  for  $i \in [k]$ :

$$\mathcal{M} = \sum_{i \in [k]} x_i \cdot \text{ct}_i = \mathbf{A} \cdot \left( \sum_{i \in [k]} x_i \mathbf{R}_i \right) + \left( \sum_{i \in [k]} x_i \mu_i \right) \cdot \mathbf{G}$$

For  $k \cdot \max_{i \in [k]} x_i \ll q$ , it can be deterministically decrypted using just the normal GSW decryption for integer message. For completeness, we briefly describe `LinearDec` as follows:

- `LinearDec(pk, sk, M)`: Compute

$$\begin{aligned} \text{sk} \cdot \mathcal{M} &= [-\bar{\mathbf{s}}^\top \mid 1] \cdot \mathbf{A} \cdot \left( \sum_{i \in [k]} x_i \mathbf{R}_i \right) + \left( \sum_{i \in [k]} x_i \mu_i \right) \cdot [-\bar{\mathbf{s}}^\top \mid 1] \cdot \mathbf{G} \\ &= \mathbf{e}^\top \cdot \left( \sum_{i \in [k]} x_i \mathbf{R}_i \right) + \left( \sum_{i \in [k]} x_i \mu_i \right) \cdot [-\bar{\mathbf{s}}^\top \mid 1] \cdot \mathbf{G} \\ &\approx \left( \sum_{i \in [k]} x_i \mu_i \right) \cdot [-\bar{\mathbf{s}}^\top \mid 1] \cdot \mathbf{G} \end{aligned}$$

Since it is easy to solve LWE with respect to  $\mathbf{G}$ , we can recover  $(\sum_{i \in [k]} x_i \mu_i) \cdot [-\bar{\mathbf{s}}^\top \mid 1]$  and hence  $\sum_{i \in [k]} x_i \mu_i$  (modulo  $q$ ).

*Efficiency Based on Subexponential LWE* We first specify the FHE parameters that we use under the subexponential LWE assumption.

Choose the lattice dimension parameter  $g = \lambda$ , the modulus  $q$  and the noise  $\sigma$  with subexponential modulus-to-noise ratio, and the sample complexity  $h \in \Theta(g \cdot \log q)$ . The subexponential hardness of LWE requires the below indistinguishability to hold for adversaries of size  $\mathcal{T}(\lambda) = 2^{\lambda^\beta}$  for some constant  $\beta > 0$ , with subexponentially small distinguishing advantage  $\epsilon(\lambda) = 2^{-\lambda^c}$  for some constant  $c > 0$  denoted by  $\approx_\epsilon$ :

$$\left\{ \begin{array}{l} \mathbf{A} \leftarrow \mathbb{Z}_q^{g \times h} \\ (\mathbf{A}, \mathbf{s} \cdot \mathbf{A} + \mathbf{e}) \pmod{q} : \mathbf{s} \leftarrow \mathbb{Z}_q^{1 \times g} \\ \mathbf{e} \leftarrow \chi_\sigma^{1 \times h} \end{array} \right\} \approx_\epsilon \left\{ (\mathbf{A}, \mathbf{u}) : \begin{array}{l} \mathbf{A} \leftarrow \mathbb{Z}_q^{g \times h} \\ \mathbf{u} \leftarrow \mathbb{Z}_q^{1 \times g} \end{array} \right\}.$$

Set the security parameter be  $\lambda = \log^r n$  for some constant  $r > 1/\beta$  and  $r > 1/c$ . We can verify this setup satisfies the relationship between functions required by the Robustness Theorem 3.4:  $\mathcal{T}(\log^r n) = 2^{(\log n)^{r\beta}} = \Omega(\text{poly } n)$  and  $\epsilon(\log^r n) = 2^{-(\log n)^{rc}} \leq \text{negl}(n)$  for some negligible functions  $\text{negl}$ .

Given these parameters, using GSW to implement PFHE, we have the following complexity:

**Proposition 4.6.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm  $\mathcal{A}$ . Assuming the above subexponential hardness of LWE, we have the following complexity regarding StreamAlg in Construction 3.1:*

1.  $(\widetilde{\text{pk}}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  and *sketch* can be stored using  $\tilde{\mathcal{O}}(1)$  total bits of space.
2.  $\mathcal{R}(\text{StreamAlg.Update}) = \tilde{\mathcal{O}}(1) + \mathcal{R}(\text{StreamAlg}_0.\text{Update})$ .
3.  $\mathcal{R}(\text{StreamAlg.Report}) = \tilde{\mathcal{O}}(k) + \mathcal{R}(\text{StreamAlg}_0.\text{Report})$ .

*Proof.* Recall that computations in StreamAlg are modulo  $q$  with  $q \in \text{poly}(n)$  and  $q \gg n \cdot N$ , where  $N \in \text{poly}(n)$  is the bound of the stream. Thus all entries have bit complexity in  $\mathcal{O}(\log n)$ .

1.  $\widetilde{\text{pk}}$ , *sketch*, and  $\widetilde{\text{ct}}_i$  for  $i \in [\log n]$  are all  $\mathcal{O}(g) \times \mathcal{O}(h)$  matrices over  $\mathbb{Z}_q$ . For  $g = \lambda \in \text{poly}(\log n)$  and  $h \in \Theta(g \cdot \log q)$ , they can be stored using  $\mathcal{O}(\text{poly}(\log n)) = \tilde{\mathcal{O}}(1)$  total bits of space.
2. The second bullet point considers the following operation:
  - Evaluate  $\widehat{\text{ct}}_{i_t} \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_{i_t}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_\ell)$ . Then update *sketch* as  $\text{sketch} \leftarrow \text{sketch} + \Delta_t \cdot \widehat{\text{ct}}_{i_t}$ .

Updating *sketch* takes  $\tilde{\mathcal{O}}(1)$  time. For evaluation, recall that  $C_{i_t}$  is a point function checking whether the input bits equal  $i_t$ . This can be naïvely implemented using  $\mathcal{O}(\log n)$  arithmetic operations, each of which runs in  $\text{poly} \log n$  time as described in  $\text{Eval}_+$  and  $\text{Eval}_\times$  above. Thus the total runtime is  $\tilde{\mathcal{O}}(1)$ .

3. The third bullet point considers the following operation:

- (If  $\|\mathbf{x}'\|_0 > k$ , output  $\perp$ . Otherwise:) Evaluate  $\widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_i, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_\ell)$  for  $i \in [n]$ . Then compute  $\text{hash} \leftarrow \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot \mathbf{x}'_i$ .

Observe that although the construction is described as evaluating  $n$ -many  $\widehat{\text{ct}}_i$ , one for each coordinate of  $\mathbf{x}'$ , at this line of the algorithm  $\mathbf{x}'$  is guaranteed to have at most  $k$  non-zero entries. Thus  $k$  evaluations (corresponding to non-zero coordinates of  $\mathbf{x}'$ ) suffice for computing the linear combination *hash*. Each evaluation is  $\tilde{\mathcal{O}}(1)$ , so the total runtime is  $\tilde{\mathcal{O}}(k)$ .

□

*Efficiency Based on Polynomial LWE* Under the polynomial LWE assumption, we set  $\lambda = n^w$  for an arbitrarily small constant  $w > 0$ . The  $(\mathcal{T}, \epsilon)$ -security then reduces to the standard definition in terms of PPT adversaries and negligible advantage. We again choose  $g = \lambda$  and  $h \in \Theta(g \cdot \log q)$ . The complexity is:

**Proposition 4.7.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm  $\mathcal{A}$ . Assuming the polynomial hardness of  $\text{LWE}$ , we have the following complexity regarding  $\text{StreamAlg}$  in Construction 3.1:*

1.  $(\widetilde{\text{pk}}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  and sketch can be stored in  $\tilde{\mathcal{O}}(n^{w_1})$  bits of space.
2.  $\mathcal{R}(\text{StreamAlg.Update}) = \tilde{\mathcal{O}}(n^{w_2}) + \mathcal{R}(\text{StreamAlg}_0.\text{Update})$ .
3.  $\mathcal{R}(\text{StreamAlg.Report}) = \tilde{\mathcal{O}}(k \cdot n^{w_3}) + \mathcal{R}(\text{StreamAlg}_0.\text{Report})$ .

where  $w_1, w_2$ , and  $w_3$  are arbitrarily small positive constants.

*Proof.* The proof is exactly the same as the proof for Proposition 4.6, except that now the dimension of the ciphertext matrices blows up from  $\text{poly log } n$  to  $n^w$  for an arbitrarily small  $w > 0$ . □

### 4.3 Complexity of Construction 3.1

Putting together Section 4.1 and 4.2, we have the following total complexity:

**Corollary 4.8.** *Assuming the subexponential hardness of  $\text{LWE}$ , the  $\text{StreamAlg}$  scheme in Construction 3.1 has the following complexity:*

- $\text{StreamAlg}$  takes  $\tilde{\mathcal{O}}(k)$  bits of space.
- $\text{StreamAlg.Update}$  runs in time  $\tilde{\mathcal{O}}(1)$ .
- $\text{StreamAlg.Report}$  runs in time  $\tilde{\mathcal{O}}(k^{1+c})$  for an arbitrarily small constant  $c > 0$ .

Corollary 4.8 follows from the efficiency of the hash-based  $\text{StreamAlg}_0$  from Theorem 4.1 and GSW. We summarize all combinations of instantiations in Table 2 below.

**Table 2.** A summary of the time and space complexities of the  $\text{StreamAlg}$  scheme in Construction 2.1, given different instantiations of  $\text{StreamAlg}_0$  and PFHE.  $w$  denotes an arbitrarily small positive constant.

REC. SCHEME	LWE ASSUMPTION	SPACE	UPDATE TIME	REPORT TIME
2-LAYER HASH.	SUBEXP	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(1)$	$\tilde{\mathcal{O}}(k^{1+w})$
VANDERMONDE	SUBEXP	$\tilde{\mathcal{O}}(k)$	$\tilde{\mathcal{O}}(1)$	$\tilde{\mathcal{O}}(k^2)$
2-LAYER HASH.	POLY	$\tilde{\mathcal{O}}(k + n^w)$	$\tilde{\mathcal{O}}(n^w)$	$\tilde{\mathcal{O}}(k \cdot (k^w + n^w))$
VANDERMONDE	POLY	$\tilde{\mathcal{O}}(k + n^w)$	$\tilde{\mathcal{O}}(n^w)$	$\tilde{\mathcal{O}}(k \cdot (k + n^w))$

## References

1. Aharonov, D., Regev, O.: Lattice problems in  $\text{np}$  comp. J. ACM **52**, 749–765 (2005)
2. Ajtai, M., Braverman, V., Jayram, T., Silwal, S., Sun, A., Woodruff, D.P., Zhou, S.: The white-box adversarial data stream model. In: Proceedings of the 41st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems. pp. 15–27. PODS ’22, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3517804.3526228>, <https://doi.org/10.1145/3517804.3526228>
3. Almeida, A., Brás, S., Sargento, S., e.a.: Time series big data: a survey on data stream frameworks, analysis and algorithms. J Big Data (2023). <https://doi.org/https://doi.org/10.1186/s40537-023-00760-1>
4. Alon, N., Ben-Eliezer, O., Dagan, Y., Moran, S., Naor, M., Yogev, E.: Adversarial laws of large numbers and optimal regret in online classification. In: STOC ’21: 53rd Annual ACM SIGACT Symposium on Theory of Computing. pp. 447–455 (2021)
5. Alon, N., Matias, Y., Szegedy, M.: The space complexity of approximating the frequency moments. In: Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing. pp. 20–29. STOC ’96, Association for Computing Machinery, New York, NY, USA (1996). <https://doi.org/10.1145/237814.237823>, <https://doi.org/10.1145/237814.237823>
6. Attias, I., Cohen, E., Shechner, M., Stemmer, U.: A framework for adversarial streaming via differential privacy and difference estimators. CoRR **abs/2107.14527** (2021)
7. Ben-Eliezer, O., Eden, T., Onak, K.: Adversarially robust streaming via dense-sparse trade-offs. In: 5th Symposium on Simplicity in Algorithms, SOSA (2022), (to appear)
8. Ben-Eliezer, O., Jayaram, R., Woodruff, D.P., Yogev, E.: A framework for adversarially robust streaming algorithms. SIGMOD Rec. **50**(1), 6–13 (2021)
9. Ben-Eliezer, O., Yogev, E.: The adversarial robustness of sampling. In: Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS. pp. 49–62 (2020)
10. Berinde, R., Indyk, P., Ruzic, M.: Practical near-optimal sparse recovery in the  $\ell_1$  norm. In: 2008 46th Annual Allerton Conference on Communication, Control, and Computing. pp. 198–205 (2008). <https://doi.org/10.1109/ALLERTON.2008.4797556>
11. Boyle, E., LaVigne, R., Vaikuntanathan, V.: Adversarially robust property-preserving hash functions. In: Blum, A. (ed.) 10th Innovations in Theoretical Computer Science Conference, ITCS 2019, January 10–12, 2019, San Diego, California, USA. LIPIcs, vol. 124, pp. 16:1–16:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2019)
12. Brakerski, Z., Gentry, C., Vaikuntanathan, V.: (leveled) fully homomorphic encryption without bootstrapping. In: Innovations in Theoretical Computer Science 2012, Cambridge, MA, USA, January 8–10, 2012. pp. 309–325 (2012)
13. Brakerski, Z., Vaikuntanathan, V.: Fully homomorphic encryption from ring-lwe and security for key dependent messages. In: Rogaway, P. (ed.) Advances in Cryptology – CRYPTO 2011. pp. 505–524. Springer Berlin Heidelberg, Berlin, Heidelberg (2011)
14. Brakerski, Z., Vaikuntanathan, V.: Circuit-ABE from LWE: Unbounded attributes and semi-adaptive security. In: Robshaw, M., Katz, J. (eds.) CRYPTO 2016,

- Part III. LNCS, vol. 9816, pp. 363–384. Springer, Heidelberg (Aug 2016). [https://doi.org/10.1007/978-3-662-53015-3\\_13](https://doi.org/10.1007/978-3-662-53015-3_13)
15. Braverman, V., Hassidim, A., Matias, Y., Schain, M., Silwal, S., Zhou, S.: Adversarial robustness of streaming algorithms through importance sampling. In: Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems, NeurIPS (2021), (to appear)
  16. Brown, P., Haas, P.J., Myllymaki, J., Pirahesh, H., Reinwald, B., Sismanis, Y.: Toward automated large-scale information integration and discovery. In: Data Management in a Connected World, Essays Dedicated to Hartmut Wedekind on the Occasion of His 70th Birthday. pp. 161–180 (2005)
  17. Chakrabarti, A., Ghosh, P., Stoeckl, M.: Adversarially robust coloring for graph streams. In: 13th Innovations in Theoretical Computer Science Conference, ITCS (2022), (to appear)
  18. Chan, T.M.: A dynamic data structure for 3-d convex hulls and 2-d nearest neighbor queries. *J. ACM* **57**(3), 16:1–16:15 (2010)
  19. Chan, T.M., He, Q.: More dynamic data structures for geometric set cover with sublinear update time. In: 37th International Symposium on Computational Geometry, SoCG. pp. 25:1–25:14 (2021)
  20. Chandrasekaran, S., Cooper, O., Deshpande, A., Franklin, M.J., Hellerstein, J.M., Hong, W., Krishnamurthy, S., Madden, S.R., Reiss, F., Shah, M.A.: Telegraphcq: Continuous dataflow processing. In: Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data. p. 668. SIGMOD ’03, Association for Computing Machinery, New York, NY, USA (2003). <https://doi.org/10.1145/872757.872857>, <https://doi.org/10.1145/872757.872857>
  21. Cormode, G., Muthukrishnan, S.: What’s new: finding significant differences in network data streams. *IEEE/ACM Transactions on Networking* **13**(6), 1219–1232 (2005). <https://doi.org/10.1109/TNET.2005.860096>
  22. Cormode, G., Garofalakis, M.: Streaming in a connected world: Querying and tracking distributed data streams. In: Proceedings of the 11th International Conference on Extending Database Technology: Advances in Database Technology. p. 745. EDBT ’08, Association for Computing Machinery, New York, NY, USA (2008). <https://doi.org/10.1145/1353343.1353442>, <https://doi.org/10.1145/1353343.1353442>
  23. Cormode, G., Jowhari, H.: Lp samplers and their applications: A survey. *ACM Comput. Surv.* **52**(1) (feb 2019). <https://doi.org/10.1145/3297715>, <https://doi.org/10.1145/3297715>
  24. Cubuk, E.D., Zoph, B., Mané, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation policies from data. *CoRR* **abs/1805.09501** (2018)
  25. Dasu, T., Johnson, T., Muthukrishnan, S., Shkapenyuk, V.: Mining database structure; or, how to build a data quality browser. In: Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data. pp. 240–251 (2002)
  26. Driscoll, J.R., Sarnak, N., Sleator, D.D., Tarjan, R.E.: Making data structures persistent. *J. Comput. Syst. Sci.* **38**(1), 86–124 (1989)
  27. Feng, Y., Woodruff, D.P.: Improved algorithms for white-box adversarial streams. In: ICML (2023)
  28. Fiat, A., Kaplan, H.: Making data structures confluent persistent. *J. Algorithms* **48**(1), 16–58 (2003)
  29. Flajolet, P., Nigel Martin, G.: Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences* **31**(2), 182–209

- (1985). [https://doi.org/https://doi.org/10.1016/0022-0000\(85\)90041-8](https://doi.org/https://doi.org/10.1016/0022-0000(85)90041-8), <https://www.sciencedirect.com/science/article/pii/0022000085900418>
30. Forbes, M.A., Shpilka, A.: On identity testing of tensors, low-rank recovery and compressed sensing. In: Proceedings of the Forty-Fourth Annual ACM Symposium on Theory of Computing. p. 163–172. STOC '12, Association for Computing Machinery, New York, NY, USA (2012). <https://doi.org/10.1145/2213977.2213995>, <https://doi.org/10.1145/2213977.2213995>
  31. Gaber, M.M.: Data Stream Processing in Sensor Networks, pp. 41–48. Springer Berlin Heidelberg, Berlin, Heidelberg (2007). [https://doi.org/10.1007/3-540-73679-4\\_4](https://doi.org/10.1007/3-540-73679-4_4), [https://doi.org/10.1007/3-540-73679-4\\_4](https://doi.org/10.1007/3-540-73679-4_4)
  32. von zur Gathen, J., Gerhard, J.: Fast polynomial evaluation and interpolation. Cambridge University Press, 3 edn. (2013). <https://doi.org/10.1017/CBO9781139856065.013>
  33. Gentry, C., Sahai, A., Waters, B.: Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based. In: Canetti, R., Garay, J.A. (eds.) CRYPTO 2013, Part I. LNCS, vol. 8042, pp. 75–92. Springer, Heidelberg (Aug 2013). [https://doi.org/10.1007/978-3-642-40041-4\\_5](https://doi.org/10.1007/978-3-642-40041-4_5)
  34. Gilbert, A.C., Li, Y., Porat, E., Strauss, M.J.: For-all sparse recovery in near-optimal time. ACM Trans. Algorithms **13**(3) (mar 2017). <https://doi.org/10.1145/3039872>, <https://doi.org/10.1145/3039872>
  35. Goldreich, O., Goldwasser, S.: On the limits of nonapproximability of lattice problems. J. Comput. Syst. Sci. **60**(3), 540–563 (2000), <https://doi.org/10.1006/jcss.1999.1686>
  36. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. CoRR **abs/1412.6572** (2014), <http://arxiv.org/abs/1412.6572>
  37. Hassidim, A., Kaplan, H., Mansour, Y., Matias, Y., Stemmer, U.: Adversarially robust streaming algorithms via differential privacy. In: Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems, NeurIPS (2020)
  38. Haviv, I., Regev, O.: Tensor-based hardness of the shortest vector problem to within almost polynomial factors. In: Johnson, D.S., Feige, U. (eds.) 39th ACM STOC. pp. 469–477. ACM Press (Jun 2007). <https://doi.org/10.1145/1250790.1250859>
  39. Holmgren, J., Liu, M., Tyner, L., Wichs, D.: Nearly optimal property preserving hashing. In: Dodis, Y., Shrimpton, T. (eds.) Advances in Cryptology - CRYPTO 2022 - 42nd Annual International Cryptology Conference, CRYPTO 2022, Santa Barbara, CA, USA, August 15-18, 2022, Proceedings, Part III. Lecture Notes in Computer Science, vol. 13509, pp. 473–502. Springer (2022)
  40. Hu, Y., Qu, A., Wang, Y., Work, D.B.: Streaming data preprocessing via online tensor recovery for large environmental sensor networks. ACM Trans. Knowl. Discov. Data **16**(6) (jul 2022). <https://doi.org/10.1145/3532189>, <https://doi.org/10.1145/3532189>
  41. Ilyas, A., Engstrom, L., Madry, A.: Prior convictions: Black-box adversarial attacks with bandits and priors. CoRR **abs/1807.07978** (2018)
  42. Indyk, P.: Stable distributions, pseudorandom generators, embeddings, and data stream computation. J. ACM **53**(3), 307–323 (2006)
  43. Indyk, P.: Explicit constructions for compressed sensing of sparse signals. In: Proceedings of the Nineteenth Annual ACM-SIAM Symposium on Discrete Algorithms. p. 30–33. SODA '08, Society for Industrial and Applied Mathematics, USA (2008)

44. Jafarpour, S.: Deterministic Compressed Sensing. Ph.D. thesis, Princeton University (2011)
45. Kacham, P., Pagh, R., Thorup, M., Woodruff, D.P.: Pseudorandom hashing for space-bounded computation with applications in streaming (2023)
46. Kane, D.M., Nelson, J., Woodruff, D.P.: On the exact space complexity of sketching and streaming small norms. In: Charikar, M. (ed.) Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin, Texas, USA, January 17-19, 2010. pp. 1161–1178. SIAM (2010)
47. Kaplan, H.: Persistent data structures. In: Handbook of Data Structures and Applications. Chapman and Hall/CRC (2004)
48. Kaplan, H., Mansour, Y., Nissim, K., Stemmer, U.: Separating adaptive streaming from oblivious streaming using the bounded storage model. In: Advances in Cryptology - CRYPTO 2021 - 41st Annual International Cryptology Conference, CRYPTO, Proceedings, Part III. pp. 94–121 (2021)
49. Khot, S.: Hardness of approximating the shortest vector problem in lattices. In: 45th FOCS. pp. 126–135. IEEE Computer Society Press (Oct 2004). <https://doi.org/10.1109/FOCS.2004.31>
50. Kurakin, A., Goodfellow, I.J., Bengio, S.: Adversarial machine learning at scale. In: 5th International Conference on Learning Representations, ICLR, Conference Track Proceedings (2017)
51. Lall, A., Sekar, V., Ogihara, M., Xu, J., Zhang, H.: Data streaming algorithms for estimating entropy of network traffic. SIGMETRICS Perform. Eval. Rev. **34**(1), 145–156 (jun 2006). <https://doi.org/10.1145/1140103.1140295>, <https://doi.org/10.1145/1140103.1140295>
52. Lightstone, S.: Physical database design for relational databases. In: Encyclopedia of Database Systems, Second Edition. Springer (2018)
53. Liu, Y., Chen, X., Liu, C., Song, D.: Delving into transferable adversarial examples and black-box attacks. In: 5th International Conference on Learning Representations, ICLR, Conference Track Proceedings (2017)
54. Lyubashevsky, V., Peikert, C., Regev, O.: On ideal lattices and learning with errors over rings. In: Gilbert, H. (ed.) EUROCRYPT 2010. LNCS, vol. 6110, pp. 1–23. Springer, Heidelberg (May / Jun 2010). [https://doi.org/10.1007/978-3-642-13190-5\\_1](https://doi.org/10.1007/978-3-642-13190-5_1)
55. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: 6th International Conference on Learning Representations, ICLR, Conference Track Proceedings (2018)
56. Menuhin, B., Naor, M.: Keep that card in mind: Card guessing with limited memory. CoRR **abs/2107.03885** (2021)
57. Muthukrishnan, S.: Data streams: Algorithms and applications. Foundations and Trends® in Theoretical Computer Science **1**(2), 117–236 (2005). <https://doi.org/10.1561/0400000002>, <http://dx.doi.org/10.1561/0400000002>
58. Nisan, N.: Pseudorandom generators for space-bounded computation. Comb. **12**(4), 449–461 (1992)
59. Phillips, J.M., Verbin, E., Zhang, Q.: Lower bounds for number-in-hand multiparty communication complexity, made easy (2015)
60. Roghani, M., Saberi, A., Wajc, D.: Beating the folklore algorithm for dynamic matching. In: 13th Innovations in Theoretical Computer Science Conference, ITCS (2022), (to appear)
61. Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., Madry, A.: Adversarially robust generalization requires more data. In: Advances in Neural Information Processing

- Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS. pp. 5019–5031 (2018)
62. Selinger, P.G., Astrahan, M.M., Chamberlin, D.D., Lorie, R.A., Price, T.G.: Access path selection in a relational database management system. In: Proceedings of the 1979 ACM SIGMOD International Conference on Management of Data. pp. 23–34. ACM (1979)
  63. Shoup, V.: Fast construction of irreducible polynomials over finite fields. *Journal of Symbolic Computation* **17**(5), 371–391 (1994). <https://doi.org/https://doi.org/10.1006/jsco.1994.1025>, <https://www.sciencedirect.com/science/article/pii/S074771718471025X>
  64. Shukla, A., Deshpande, P., Naughton, J.F., Ramasamy, K.: Storage estimation for multidimensional aggregates in the presence of hierarchies. In: VLDB’96, Proceedings of 22th International Conference on Very Large Data Bases. pp. 522–531 (1996)
  65. Smart, N., Vercauteren, F.: Fully homomorphic SIMD operations. *Cryptology ePrint Archive, Report 2011/133* (2011), <https://eprint.iacr.org/2011/133>
  66. Sung, M., Kumar, A., Li, L., Wang, J., Xu, J.: Scalable and efficient data streaming algorithms for detecting common content in internet traffic. In: 22nd International Conference on Data Engineering Workshops (ICDEW’06). pp. 27–27 (2006). <https://doi.org/10.1109/ICDEW.2006.130>
  67. Szalay, A., Gray, J., VandenBerg, J.: Petabyte scale data mining: Dream or reality? *Proceedings of SPIE - The International Society for Optical Engineering* (08 2002). <https://doi.org/10.1117/12.461427>
  68. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I.J., Fergus, R.: Intriguing properties of neural networks. *International Conference on Learning Representations* (2014)
  69. Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I.J., Boneh, D., McDaniel, P.D.: Ensemble adversarial training: Attacks and defenses. In: 6th International Conference on Learning Representations, ICLR, Conference Track Proceedings (2018)
  70. Venkataraman, S., Song, D.X., Gibbons, P.B., Blum, A.: New streaming algorithms for fast detection of superspreaders. In: *Network and Distributed System Security Symposium* (2005), <https://api.semanticscholar.org/CorpusID:538509>
  71. Wajc, D.: Rounding dynamic matchings against an adaptive adversary. In: *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC*. pp. 194–207 (2020)
  72. Woodruff, D.P., Zhang, Q.: Tight bounds for distributed functional monitoring. In: *Proceedings of the Forty-Fourth Annual ACM Symposium on Theory of Computing*. p. 941–960. STOC ’12, Association for Computing Machinery, New York, NY, USA (2012). <https://doi.org/10.1145/2213977.2214063>, <https://doi.org/10.1145/2213977.2214063>
  73. Woodruff, D.P., Zhou, S.: Tight bounds for adversarially robust streams and sliding windows via difference estimators. In: *62nd IEEE Annual Symposium on Foundations of Computer Science, FOCS*. pp. 1183–1196 (2021)
  74. Zhao, Q.G., Kumar, A., Wang, J., Xu, J.J.: Data streaming algorithms for accurate and efficient measurement of traffic and flow matrices. *SIGMETRICS Perform. Eval. Rev.* **33**(1), 350–361 (jun 2005). <https://doi.org/10.1145/1071690.1064258>, <https://doi.org/10.1145/1071690.1064258>



## Supplementary Material

We discuss our results in the distributed model, and our results for matrix and tensor recovery in the supplementary material.

## A Distributed $K$ -Sparse Recovery from Ring LWE

### A.1 Distributed Computation

In the distributed computation model, the dataset is split over multiple servers. The goal is to design a communication protocol for servers to collectively compute a function on the aggregated dataset. In this work, we consider the message-passing coordinator model (see, e.g., [59]), where there are  $T$  servers and a single central coordinator. Each server  $s_t : t \in [T]$  receives a partition of data in the form of an  $n$ -dimensional vector  $\mathbf{x}_t$ , and communicates with the central coordinator through a two-way, single-round communication channel. At the end of communication, the coordinator needs to compute or approximate a function  $f_n(\mathbf{x})$  of the aggregated vector  $\mathbf{x} = \sum_{t \in [T]} \mathbf{x}_t$ . We aim to design protocols that minimize the total amount of communication through the channel and speed up the computation at each server and the coordinator. See also [72] and references therein on the distributed model.

Similar to the streaming case, we focus on designing a robust protocol in a white-box adversarial distributed setting, where the dataset is generated and partitioned to each server by a white-box adversary who monitors the internal state of the coordinator and the servers. We define the model as follows:

**Definition A.1 (Distributed Protocol).** *Given a query function ensemble  $\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow X\}_{n \in \mathbb{N}}$  for some domain  $X$ , a distributed protocol  $\text{DistProt}$  that computes  $\mathcal{F}$  contains a tuple of algorithms  $\text{DistProt} = (\text{DistProt.Setup}, \text{DistProt.Server}, \text{DistProt.Coordinator})$ , with the following properties:*

- $\text{DistProt.Setup}(n, \rho_0) \rightarrow \zeta$  : *On input a parameter  $n$ , the coordinator samples randomness  $\rho_0$  and uses  $\rho_0$  to compute some state  $\zeta$ . It sends  $\zeta$  to all servers.*
- $\text{DistProt.Server}(\zeta, \mathbf{x}_t, \rho_t) \rightarrow p_t$  : *On input a state  $\zeta$  and a data partition  $\mathbf{x}_t \in \mathbb{Z}^n$ , the  $t$ -th server samples randomness  $\rho_t$  and uses it to compute a packet  $p_t$ . It sends  $p_t$  back to the coordinator.*
- $\text{DistProt.Coordinator}(p_1, \dots, p_T, \rho_c) \rightarrow r$  : *On input  $T$  packets  $p_1, \dots, p_T$ , the coordinator samples randomness  $\rho_c$  and computes a query response  $r \in X$ . A response  $r$  is said to be correct if  $r = f_n(\mathbf{x})$  where  $\mathbf{x} = \sum_{t \in [T]} \mathbf{x}_t$ .*

As in the streaming setting, we again assume bounded integer inputs, such that an arbitrary subset sum of  $\{\mathbf{x}_t : t \in [T]\}$  is in  $[-N, N]^n$  for some  $N \in \text{poly}(n)$ .

**Definition A.2.** *We say that a distributed protocol  $\text{DistProt}$  is white-box adversarially robust if for all dimension parameters  $n \in \mathbb{N}$ , the following holds: For*

any PPT adversary  $\mathcal{A}$ , the following experiment  $\text{Expt}_{\mathcal{A}, \text{DistProt}}(n)$  outputs 1 with probability at most  $\text{negl}(n)$ :

$\text{Expt}_{\mathcal{A}, \text{DistProt}}(n)$  :

1. The challenger and  $\mathcal{A}$  agree on a query function ensemble  $\mathcal{F} = \{f_n : \mathbb{Z}^n \rightarrow X\}_{n \in \mathbb{N}}$  for some domain  $X$ . On input a dimension parameter  $n$ , let  $f_n \in \mathcal{F}$  be the query function of the experiment.
2. The challenger samples all randomness  $\rho_0, \rho_t : t \in [T]$ , and  $\rho_c$  for the servers and the coordinator. It then provides  $\rho_0, \rho_1, \dots, \rho_T$  and  $\rho_c$  to  $\mathcal{A}$ .
3.  $\mathcal{A}$  generates  $T$  partitions of data  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{Z}^n$  and provides them to the challenger.
4. The challenger executes the protocol to sequentially compute the following:
  - (i)  $\zeta \leftarrow \text{DistProt.Setup}(n, \rho_0)$ .
  - (ii)  $p_t \leftarrow \text{DistProt.Server}(\zeta, \mathbf{x}_t, \rho_t)$  for all  $t \in [T]$ .
  - (iii)  $r \leftarrow \text{DistProt.Coordinator}(p_1, \dots, p_T, \rho_c)$ .
5. The experiment outputs 1 if and only if at some time  $r \neq f_n(\mathbf{x})$ .

*Remark A.3.* We remark on the connection between the streaming model and the distributed model. There exists a well-known reduction from a streaming algorithm to a distributed protocol under the oblivious setting: Each server runs  $\text{StreamAlg.Update}$  on the same initial state  $\zeta \leftarrow \text{StreamAlg.Setup}$ , using all entries of  $\mathbf{x}_t$  (its data partition) as updates. It then sends the resulting sketch to the coordinator. As long as the sketches are linearly measured, the coordinator can run  $\text{StreamAlg.Report}$  on the sum of all sketches to compute a query response.

One may use the same recipe to construct a *WAR* distributed protocol from a *WAR* streaming algorithm. In general, the robustness is retained by such reduction if  $\text{StreamAlg.Update}$  and  $\text{StreamAlg.Report}$  are deterministic. (Otherwise, adversaries in the experiment  $\text{Expt}_{\mathcal{A}, \text{DistProt}}$  would have additional views to all “future” randomness, while adversaries in the streaming experiment  $\text{Expt}_{\mathcal{A}, \text{StreamAlg}}$  only monitors the past randomness already used by  $\text{StreamAlg}$ .) Our Construction 3.1 satisfies this condition.

## A.2 Distributed Protocol Construction

For completeness, we briefly describe the distributed protocol for  $k$ -Sparse Recovery as follows, which essentially follows the same route as our streaming Construction 3.1.

**Construction A.1.** *Similar to Construction 3.1, let PFHE be a pseudorandom FHE scheme. And let  $\text{StreamAlg}_0$  be a streaming algorithm for the Relaxed  $k$ -Sparse Recovery problem. We assume that  $\text{StreamAlg}_0$  performs a deterministic linear measurement  $\Phi_n$  on arbitrary length- $n$  input vector, as described in the previous instantiation Section 4.1.*

- $\text{DistProt.Setup}(n)$  : Let  $\lambda = \mathcal{S}(n)$  be the security parameter of PFHE. Sample  $\widetilde{\text{pk}} \leftarrow \mathcal{D}_{\widetilde{\text{pk}}(\lambda)}$  and  $\widetilde{\text{ct}}_1, \widetilde{\text{ct}}_2, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil} \leftarrow \mathcal{D}_{\widetilde{\text{ct}}(\lambda)}$ . Send  $\zeta := (\widetilde{\text{pk}}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  to all servers

- $\text{DistProt.Server}(\zeta, \mathbf{x}_t) : \text{Given } \zeta := (\widetilde{\text{pk}}, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil}), \text{ evaluate } \widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_i, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil}) \text{ for } i \in [n], \text{ where } C_i \text{ is the same point function as in Construction 3.1.}$   
 $\text{Compute } \Phi_n(\mathbf{x}_t) \text{ and } \text{sketch}_t \leftarrow \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot (\mathbf{x}_t)_i, \text{ where } (\mathbf{x}_t)_i \text{ denotes the } i\text{-th coordinate of } \mathbf{x}_t. \text{ Send } p_t := (\text{sketch}_t, \Phi_n(\mathbf{x}_t)) \text{ to the coordinator.}$
- $\text{DistProt.Coordinator}(p_1, \dots, p_T) : \text{Let } \mathbf{x}' \leftarrow \text{StreamAlg}_0.\text{Report}(\sum_{t \in [T]} \Phi_n(\mathbf{x}_t)), \text{ i.e. the output of running the reconstruction algorithm of } \text{StreamAlg}_0 \text{ on the sum of all measurement results. If } \|\mathbf{x}'\|_0 > k, \text{ output } \perp.$   
 $\text{Otherwise, evaluate } \widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_i, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil}) \text{ for } i \in [n].$   
 $\text{Then compute } \text{hash} \leftarrow \sum_{i \in [n]} \widehat{\text{ct}}_i \cdot \mathbf{x}'_i. \text{ Output } \mathbf{x}' \text{ if } \text{hash} = \sum_{t \in [T]} \text{sketch}_t, \text{ and output } \perp \text{ otherwise.}$

### A.3 Efficiency

In this section, we focus on showing that we can use BV [13], a Ring-LWE based FHE scheme, to achieve near-linear process time on  $\text{DistProt.Server}$  under *polynomial ring-LWE* assumption.

As a warm-up, we state the complexity of Construction A.1 using GSW to implement PFHE:

**Proposition A.4.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm or an evaluation  $\mathcal{A}$ , we have the following complexity regarding  $\text{DistProt}$  in Construction A.1:*

- *Assuming the subexponential hardness of LWE:*
  1. *The total bits of communication between an arbitrary server and the coordinator is  $\tilde{\mathcal{O}}(1) + \text{size}(\Phi_n(\mathbf{x}))$ .*
  2.  $\mathcal{R}(\text{DistProt.Coordinator}) = \tilde{\mathcal{O}}(k) + \mathcal{R}(\text{StreamAlg}_0.\text{Report})$ .
  3.  $\mathcal{R}(\text{DistProt.Server}) = \tilde{\mathcal{O}}(k) + \mathcal{R}(\Phi_n(\mathbf{x}))$ .
- *Assuming the polynomial hardness of LWE:*
  1. *The total bits of communication between an arbitrary server and the coordinator is  $\tilde{\mathcal{O}}(n^{w_1}) + \text{size}(\Phi_n(\mathbf{x}))$ .*
  2.  $\mathcal{R}(\text{DistProt.Coordinator}) = \tilde{\mathcal{O}}(k \cdot n^{w_2}) + \mathcal{R}(\text{StreamAlg}_0.\text{Report})$ .
  3.  $\mathcal{R}(\text{DistProt.Server}) = \tilde{\mathcal{O}}(k \cdot n^{w_3}) + \mathcal{R}(\Phi_n(\mathbf{x}))$ .

where  $\mathbf{x} \in [-N, N]^n$  is arbitrary, and  $w_1, w_2$ , and  $w_3$  are arbitrarily small positive constants.

*Proof.* The communication complexity and the runtime of  $\text{DistProt.Coordinator}$  are the same as the bit complexity and report time of  $\text{StreamAlg}$  in Construction 3.1 (except some negligible addition operations). The runtime of  $\text{DistProt.Server}$  is roughly  $n \cdot \mathcal{R}(\text{StreamAlg.Update})$ , because every server has to perform  $n$  evaluations, one for each coordinate of its data partition.  $\square$

*Efficiency Based on Ring-LWE.* Before discussing the efficiency advantage, we briefly review some important features of BV to show that it satisfies a variant of the PFHE definition. See [13] for details of the BV construction.

Consider a ring  $R := \frac{\mathbb{Z}[x]}{x^g+1}$  where  $g$  is the degree parameter. Let  $R_q$  denote  $R/qR$  for a prime modulus  $q$ , and let  $R_t$  be the message space for some prime  $t \in \mathbb{Z}_q^*$  and  $t \ll q/g$ :

- The public keys of BV are two-dimensional vectors in  $R_q^2$ , which are computationally indistinguishable from uniformly random samples from  $R_q^2$  under the Ring-LWE assumption.  
The ciphertexts are also vectors over  $R_q$  that are indistinguishable from uniform random.
- BV satisfies an analogy of  $\mathbb{Z}$  Linear Homomorphism, with coefficients of linear combinations  $\mathcal{M} = \sum_{i \in [k]} x_i \cdot \text{ct}_i$  being small norm ring elements  $x_i \in R_t$  as opposed to integer scalars. As long as the linear combination of corresponding messages ( $\mathcal{M}_{\text{msg}} = \sum_{i \in [k]} x_i \cdot m_i$ , where  $m_i \in R_t$  is the message encrypted by  $\text{ct}_i$ ) still lies in  $R_t$ , running the standard BV decoding algorithm on  $\mathcal{M}$  output the correct message.
- Eval and Dec are deterministic by the construction of BV.

We choose parameters for our Ring-LWE based construction as follows: Let  $\lambda = n^w$  for an arbitrarily small constant  $w > 0$ . Let the dimension  $g = \lambda$ , the modulus of the message space  $t > n \cdot N$ , where  $N \in \text{poly}(n)$  is the bound of the input data. The modulus  $q$  is chosen with subexponential modulus-to-noise ratio and  $q \gg g \cdot t$ . The Ring-LWE assumption states that for any  $h \in \text{poly}(\lambda)$ , it holds that

$$\{(a_i, a_i \cdot s + e_i) : e_i \leftarrow \chi\}_{i \in [h]} \approx_c \{(a_i, u_i) : u_i \leftarrow R_q\}_{i \in [h]}$$

where  $s$  is sampled from the noise distribution  $\chi$  and  $a_i$  are uniform in  $R_q$ . We have the following complexity:

**Proposition A.5.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm or an evaluation  $\mathcal{A}$ . Assuming the polynomial hardness of Ring-LWE, a variant of DistProt in Construction A.1 has the following complexity:*

1. *The total bits of communication between an arbitrary server and the coordinator is  $\tilde{O}(n^{w_1}) + \text{size}(\Phi_n(\mathbf{x}))$ .*
2.  $\mathcal{R}(\text{DistProt.Coordinator}) = \tilde{O}(n) + \mathcal{R}(\text{StreamAlg}_0.\text{Report})$ .
3.  $\mathcal{R}(\text{DistProt.Server}) = \tilde{O}(k \cdot n_2^w) + \mathcal{R}(\Phi_n(\mathbf{x}))$ .

where  $\mathbf{x} \in [-N, N]^n$  is arbitrary, and  $w_1, w_2$  are arbitrarily small positive constants.

*Proof.* In comparison to the complexity under polytime standard LWE, we save a  $n^{\Theta(1)}$  factor in the runtime of DistProt.Server. This is due to speeding up the following two steps:

1. Evaluate  $\widehat{\text{ct}}_i \leftarrow \text{PFHE.Eval}(\widetilde{\text{pk}}, C_i, \widetilde{\text{ct}}_1, \dots, \widetilde{\text{ct}}_{\lceil \log n \rceil})$  for all  $i \in [n]$

2. Compute  $\sum_{i \in [n]} \widehat{\mathbf{ct}}_i \cdot \mathbf{x}_i$

The idea is to represent the messages as ring elements over  $R_q$  as opposed to individual scalars, and apply fast Fourier transforms to efficiently multiply them.

In a variant of Construction A.1, instead of performing  $n$  evaluations one for each coordinate  $i \in [n]$ , we partition the dimension  $n$  into  $n^{1-w}$  chunks. Each chunk is of length  $n^w$ , where  $g = n^w$  is the degree parameter of the ring for an arbitrarily small constant  $w > 0$ . Such that a length- $n$  vector  $\mathbf{x}$  is partitioned into  $n^{1-w}$  consecutive segments, each can be represented as a (small-normed in comparison to  $q$ ) ring element  $x_j \in R_q$  for  $j \in [n^{1-w}]$ .

We also define a new set of functions  $C'_j$  as follows:

$$C'_j(\mu_1, \dots, \mu_{\lceil \log n \rceil}) = \begin{cases} 1 & \text{if } (\mu_1, \dots, \mu_{\lceil \log n \rceil}) \text{ represents a value in } [jn^w, (j+1)n^w) \\ 0 & \text{otherwise.} \end{cases}$$

We evaluate  $\widehat{\mathbf{ct}}'_j \leftarrow \text{PFHE.Eval}(\widetilde{\mathbf{pk}}, C'_j, \widetilde{\mathbf{ct}}_1, \dots, \widetilde{\mathbf{ct}}_{\lceil \log n \rceil})$  for all  $j \in [n^{1-w}]$  to replace the above step 1. The total evaluation time is  $\mathcal{O}(n^{1-w} \cdot n^w \text{poly log}(n^w)) = \tilde{\mathcal{O}}(n)$ , since two ring elements can be multiplied in  $\mathcal{O}(n^w \text{poly log}(n^w))$  time using FFT. Then the linear combination in the above step 2 is replaced by  $\sum_{j \in [n^{1-w}]} \widehat{\mathbf{ct}}'_j \cdot x_j$ , which similarly takes  $\tilde{\mathcal{O}}(n)$  time to compute.

Lastly, we check that given two vectors  $\mathbf{x}$  and  $\mathbf{x}'$  differ at a coordinate  $m$ , when  $\mathbf{ct}_1, \dots, \mathbf{ct}_{\lceil \log n \rceil}$  encrypt the bits of  $m$ , we still have

$$\text{sketch} = \sum_{j \in [n^{1-w}]} \widehat{\mathbf{ct}}'_j \cdot x_j \neq \sum_{j \in [n^{1-w}]} \widehat{\mathbf{ct}}'_j \cdot x'_j = \text{hash}.$$

because they can be decrypted into distinct ring elements. This concludes that the above variant of Construction A.1 remains to be correct while achieving near-linear runtime of `DistProt.Coordinator` and `DistProt.Server`.  $\square$

## B Low-Rank Matrix and Tensor Recovery

The low-rank matrix and tensor recovery problems are defined as follows:

**Definition B.1 (Rank- $k$  Matrix Recovery).** *Given a rank parameter  $k$ , let  $\mathbf{X} \in \mathbb{Z}^{n \times m}$  denote an input data matrix. The query function ensemble of the Rank- $k$  Matrix Recovery problem is:*

$$\mathcal{F} = \{f_n : \mathbb{Z}^{n \times m} \rightarrow (\mathbb{Z}^{n \times m} \cup \{\perp\})\}_{n \geq m \in \mathbb{N}},$$

$$\text{where } f_n(\mathbf{X}) = \begin{cases} \mathbf{X} & \text{if } \text{rank}(\mathbf{X}) \leq k \\ \perp & \text{otherwise} \end{cases}$$

Let  $\otimes$  denote the outer product of two vectors. We use the following definition for tensor rank:

**Definition B.2 (CANDECOMP/PARAFAC (CP) rank).** *For a tensor  $X \in \mathbb{R}^{n^d}$ , consider it to be the sum of  $k$  rank-1 tensors:  $X = \sum_{i=1}^k (x_{i1} \otimes x_{i2} \otimes \cdots \otimes x_{id})$  where  $x_{ij} \in \mathbb{R}^n$ . The smallest number of rank-1 tensors that can be used to express a tensor  $X$  is then defined to be the rank of the tensor.*

**Definition B.3 (Rank- $k$  Tensor Recovery).** *Given a rank parameter  $k$ , let  $\mathbf{X} \in \mathbb{Z}^{n^d}$  denote an input data tensor. The query function ensemble of the Rank- $k$  Tensor Recovery problem is:*

$$\mathcal{F} = \{f_n : \mathbb{Z}^{n^d} \rightarrow (\mathbb{Z}^{n^d} \cup \{\perp\})\}_{n \in \mathbb{N}, d \in \mathcal{O}(1)},$$

$$\text{where } f_n(\mathbf{X}) = \begin{cases} \mathbf{X} & \text{if } \text{rank}(\mathbf{X}) \leq k \\ \perp & \text{otherwise} \end{cases}$$

The low-rank matrix and tensor recovery problems are natural generalizations of the  $k$ -sparse vector recovery problems. Using the same framework as the vector case, they can be solved by maintaining two hashes: one recovers the matrix or tensor but potentially produces garbage output, and the other uses a cryptographic tool to verify whether the output matches the input data. These are natural generalizations of the  $k$ -sparse vector recovery problems.

In [27], the recovery hash is measured by a random Gaussian matrix, and the verification hash is measured by an SIS matrix. Both measurements are assumed to be heuristically compressed by an idealized hash function. We instead replace the former with a deterministic measurement, which can be explicitly constructed and decoded in polynomial time [30]. The latter is replaced by the same FHE-based hashing as the one we used for the vector case in Construction 3.1, which operates on the vectorization of the data matrix or tensor. Similar to the vector case, our construction does not require using an idealized hash function or storing a separate large state.

In the streaming model, the complexities of our matrix and tensor recovery algorithms are as follows:

**Proposition B.4.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm or an evaluation  $\mathcal{A}$ . Assuming the polynomial hardness of LWE,*

- *There exists a WAR streaming algorithm **StreamAlg** for the rank- $k$  matrix recovery problem using  $\tilde{\mathcal{O}}(nk)$  bits of space.*
- *There exists a WAR streaming algorithm **StreamAlg** for the rank- $k$  tensor recovery problem using  $\tilde{\mathcal{O}}(nk^{\lceil \log d \rceil})$  bits of space.*

*Proof.* The recovery scheme for matrices and tensors takes  $\tilde{\mathcal{O}}(nk)$  and  $\tilde{\mathcal{O}}(nk^{\lceil \log n \rceil})$  bits of space, respectively [30]. Under the polynomial hardness assumption of LWE, we set  $\lambda = n^c$  for an arbitrarily small positive constant  $c$ . Taking the sum of the two terms produces the claimed complexities.  $\square$

In addition, the above streaming algorithms can be converted into WAR distributed protocols for matrix and vector recovery, following the same roadmap as in Section A.

**Proposition B.5.** *Let  $\mathcal{R}(\mathcal{A})$  denote the runtime of an algorithm or an evaluation  $\mathcal{A}$ . Assuming the polynomial hardness of LWE: :*

- *There exists a WAR distributed protocol **DistProt** for the rank- $k$  matrix recovery problem, using  $\tilde{\mathcal{O}}(nk)$  bits of communication between each server and the coordinator, and  $\tilde{\mathcal{O}}(n^2k)$  processing time at each server.*
- *There exists a WAR streaming algorithm **StreamAlg** for the rank- $k$  tensor recovery problem, using  $\tilde{\mathcal{O}}(nk^{\lceil \log d \rceil})$  bits of communication between each server and the coordinator.*

*Proof.* The total bits of communication between each server and the coordinator equals the size of the state in the above streaming algorithms. For the low-rank matrix recovery, the deterministic scheme performs  $\tilde{\mathcal{O}}(nk)$  measurements using  $n$ -sparse matrices, and the runtime for computing the FHE-based hash is comparably small. This gives  $\mathcal{O}(n^2k)$  processing time at each server.  $\square$

We remark that the rank- $k$  matrix recovery algorithm in the previous work [27] requires  $\tilde{\mathcal{O}}(n^3k)$  processing time in the distributed setting, which our construction improves upon.