

# SANNS: Scaling Up Secure Approximate $k$ -Nearest Neighbors Search

Hao Chen    Iliaria Chillotti    Yihe Dong    Oxana Poburinnaya    Ilya Razenshteyn    M. Sadegh Riazi  
Microsoft Research    KU Leuven    Microsoft Research    Boston University    Microsoft Research    UC San Diego

**Abstract**—We present new secure protocols for approximate  $k$ -nearest neighbor search ( $k$ -NNS) over the Euclidean distance in the semi-honest model. Our implementation is able to handle massive datasets efficiently. On the algorithmic front, we show a new circuit for the approximate top- $k$  selection from  $n$  numbers that is built from merely  $O(n + \text{poly}(k))$  comparators. Using this circuit as a subroutine, we design new approximate  $k$ -NNS algorithms and two corresponding secure protocols: 1) optimized linear scan; 2) clustering-based *sublinear time* algorithm.

Our secure protocols utilize a combination of additively-homomorphic encryption, garbled circuit and Oblivious RAM. Along the way, we introduce various optimizations to these primitives, which drastically improve concrete efficiency.

We evaluate the new protocols empirically and show that they are able to handle datasets that are significantly larger than in the prior work. For instance, running on two standard Azure instances within the same availability zone, for a dataset of 96-dimensional descriptors of 10 000 000 images, we can find 10 nearest neighbors with average accuracy 0.9 in under 10 seconds improving upon prior work by at least two orders of magnitude.

## I. INTRODUCTION

The  $k$ -Nearest Neighbor Search problem ( $k$ -NNS from now on) can be defined as follows. For a given dataset  $X \subset \mathbb{R}^d$  lying in a  $d$ -dimensional space, and a query point  $q \in \mathbb{R}^d$ , the goal is to find  $k$  data points closest (with respect to the Euclidean distance) to the query. To improve the search efficiency, one typically relaxes the  $k$ -NNS problem in two ways. First, one allows the answer to be *approximate* (i.e., the returned set of  $k$  points should contain most but not necessarily all of the true  $k$  closest points). Second, one may allow the one-time *preprocessing phase*, during which some auxiliary information is computed, which can be later used to speed up the query procedure.

The  $k$ -NNS has many applications in modern data analysis, including web search, face recognition, recommendation systems, advertisement matching, drug design, DNA analysis, plagiarism detection, motion planning, spell checking, machine learning and other areas. One typically starts with a dataset and, using domain expertise together with machine learning tools, produces *feature vector representation* of the dataset. Then, *similarity search* queries (“find  $k$  objects most similar to the query”) directly translate to  $k$ -NNS queries in the feature space. Let us note, as a side remark, that one standard modern technique for producing feature vectors is to train a deep neural network and then read off the feature values from one of the layers [1].

When it comes to applications dealing with sensitive data, such as medical, biological or financial data, the privacy of the

information contained in the dataset and the queries needs to be ensured. One can naturally pose the *Secure  $k$ -NNS* problem, where a *server* stores the dataset  $X \subset \mathbb{R}^d$ , and a *client* holds one or several query points  $q \in \mathbb{R}^d$ . We would like the client to learn  $k$  data points (approximately) closest to  $q$  such that the server learns nothing about the query or the result, while the client should not learn anything about the dataset besides the answer to the query.

From the theoretical viewpoint problems like these are well-understood: one can use secure two-party computation protocols [2], [3] or homomorphic encryption [4], [5], [6], [7], [8], [9]. However, the known generic constructions of these primitives as of today do not lead to practically efficient protocols for the Secure  $k$ -NNS problem. As a result, secure  $k$ -NNS has been thoroughly studied on its own: see Section I-B for an overview.

In this paper, we design and evaluate two new highly-efficient and secure  $k$ -NNS protocols. The first protocol is a secure implementation of the (heavily-optimized) *linear scan*, where we compute distances from the query to all the data points, and then choose  $k$  smallest ones. The second protocol is based on a new *sublinear-time*  $k$ -NNS algorithm, which avoids computing all the distances. The new algorithm is based on *clustering*: at a very high level, during the preprocessing phase, we cluster the dataset, and then during the query stage, we search for closest points in several clusters that are the closest to the query point.

**Security guarantee.** The security of approximate  $k$ -NNS can be defined in several ways. In this work, we follow the standard approach and require that the secure protocol does not reveal more than what is revealed by the *outputs* of a *plaintext* approximate  $k$ -NNS algorithm<sup>1</sup>. Note that approximate answer by itself can potentially reveal more than what the exact answer would.

We remark that in the clustering-based protocol, the client does learn the *hyperparameters*: for example, the total number of clusters, or the number of clusters the protocol processes during the query stage (see Section III-F for more details). Even though the hyperparameters can a priori be arbitrary, the client can expect the server to set them in a way that optimizes the performance of the overall computation. We leave to future work the task of analyzing the potential leakage from these hyperparameters or hiding them from the client. Note that this situation is similar to the line of work on secure inference of

<sup>1</sup>For an alternative definition of security in this setting see [10].

neural networks (e.g. [11], [12]), where the hyperparameters of the underlying neural network, such as number of layers and number of nodes in each layer, are revealed to both parties.

Our construction can be proven secure in the semi-honest model, where both parties follow the protocol honestly.

#### A. Our contributions

*a) Plaintext approximate  $k$ -NNS algorithms tailored to secure computation:* There is a huge body of work on  $k$ -NNS algorithms (both theoretical and practical): see [13], [14], [15] for an overview of the area. However, those protocols are not tailored to be efficient in the context of secure computation. For instance, consider the task of hiding the database access pattern, which is necessary to prevent the server learning information about the query. In algorithms which access all data points in a *coherent* way – e.g. in those which do a linear scan – this is not an issue; however, (non-secure) algorithms that currently perform the best [16] are not scanning the entire dataset, and therefore one would have to employ Oblivious RAM (ORAM) to hide the access pattern. The issue is the best-performing algorithms, which are based on following paths in certain carefully constructed graphs, are highly adaptive. Hence, when implemented in secure computation, they would require many rounds of interaction, each protected by ORAM, which makes them inefficient. Another issue which greatly affects performance is that the algorithm from [16] and related ones are not "regular": that is, they adaptively compute many *individual* distances, rather than doing the same computation on large batches of points in the dataset. This does not leave any room for certain optimizations which would be possible otherwise, e.g., batching the computation using vectorization in homomorphic encryption.

These observations give us two natural ways for solving our problem. Our first algorithm does a linear scan of the database to compute all distances to the query point and returns the  $k$  closest points. To achieve good performance, we employ a number of algorithmic and implementational optimizations. In particular, we introduce an efficient circuit that performs approximate top- $k$  selection which greatly impacts the overall search performance.

Our second algorithm is *sublinear*, and it is specifically designed to perform relatively few non-adaptive memory accesses and compute distances to many points at once. Our starting point is a classic *clustering-based* approach, which appears in [17] and relies on the *k-means* clustering of the dataset. In short, during the query stage, we find several clusters that are closest to the query, and choose closest points from these clusters as an answer. The problem with this algorithm is that resulting clusters are highly unbalanced in cardinality, which adversely affects performance, since we would have to pad all clusters to the same size to avoid information leakage. In order to rectify this issue, we perform clustering iteratively at different scales, making sure the resulting clusters are balanced in size and changing the query procedure accordingly. See Section III-B for more details.

Finally, let us note that the set of primitives developed in this paper should be sufficient to implement many other  $k$ -NNS algorithms such as locality-sensitive hashing (LSH) [18]. We plan to investigate this direction in the future work.

*b) Approximate top- $k$  selection:* Both of our algorithms rely extensively on the top- $k$  selection: given a (secret-shared) sequence of  $n$  numbers of  $b$  bits each, find  $k$  smallest of them. In order to implement this in a secure way, we need to design a Boolean *circuit* that performs the top- $k$  selection. If  $k = 1$ , this is easily done in optimal  $O(bn)$  gates, since we just need to compute the minimum, however, the question for  $k > 1$  becomes more interesting. In all the prior work on secure  $k$ -NNS either only the case  $k = 1$  was considered, or the naïve circuits of sizes  $O(bnk)$  or  $O(b^2n)$  have been used. One can use sorting networks and obtain the bound of  $O(bn \log k)$  gates.

In this work, we show a new *randomized* circuit for top- $k$  selection with only  $O(b \cdot (n + \text{poly}(k)))$  gates, which outputs the correct result with high probability. The circuit is extremely simple and practical, and gives a large boost in the overall performance. As a result, even our implementation of the linear scan already significantly improves upon the prior work for, say,  $k = 10$ . We give theoretical analysis of the accuracy of the circuit which we confirm with the numerical simulations. We also note that special precautions must be taken to enforce the security of the resulting protocol that uses the new circuit, since it is *randomized* (See Section IV).

*c) Hybrid secure protocol:* Both of our algorithms comprise of two major subroutines: computing distances between a query and a list of points, and top- $k$  selection. In case of clustering-based algorithm, we also require random memory accesses. We implement distance computation using additively-homomorphic encryption (AHE), top- $k$  selection using garbled circuits (GC) and random access via distributed oblivious RAM (DORAM). Removing any of these primitives results in significantly inferior overall performance. For AHE, we use the SEAL library [19], which implements the BFV scheme [7], for garbled circuits we use our own implementation of Yao's protocol [2], and for DORAM we implement the Floram construction in read-only mode [20].

*d) Optimizing the cryptographic primitives:* We made special-purpose optimizations to the underlying cryptographic primitives to improve efficiency of our protocol. Most notably, in the Floram construction we replace AES with Kreyvium [21], which allows us to reduce the communication and computation of DORAM by more than an order of magnitude. When we use AHE for computing distances to a list of points, we utilize the SIMD capabilities of the BFV scheme. Our approach allows us to avoid expensive rotation operations and, at the same time, set the plaintext modulus to a power of two. The latter makes the top- $k$  part of the algorithm substantially faster.

*e) Efficient implementation:* We implement our protocols in 7400 lines of C++ code and evaluate them on two datasets: SIFT [22] (1 000 000 image descriptors) and more modern Deep1B [23] (1 000 000 000 image descriptors obtained using

deep neural networks, from which we subsample 1 and 10 million). We require to return 10 nearest neighbors so that on average 9 of them are correct (accuracy 0.9). We find that the clustering-based algorithm is faster than the linear scan, on the largest dataset by more than an order of magnitude. Yet linear scan itself is faster than the prior work by at least an order of magnitude due to a better top- $k$  circuit.

Overall, we show the first practically efficient secure implementation of a *sublinear-time* NNS algorithm, and our work is the first to handle datasets of the scale of tens of millions points, whereas we are not aware of any prior work, which runs secure NNS on datasets more than several thousand points.

## B. Related work

The works [24], [25], [26], [27], [28] consider the secure computation scenarios that can be mapped to the  $k$ -NNS problem for  $k = 1$ , with the exception that [27] returns all matches with distance below a given threshold. While these works employ different techniques, they share some common properties: first, they perform linear scan over the database. Second, these works use the Paillier AHE scheme [29] for distance computation (except for [28], which uses secret sharing schemes). In contrast, we use a more recently developed packed lattice-based AHE scheme which significantly reduces the computation cost. Moreover, all experiments done in these works have the servers database size to be at most 5 000.

Several works implemented secure algorithms tailored for NNS. The work [30] assumes that both dataset and query belong to the client and the goal is to outsource the search computation to a server. However, this results in significant drawbacks in efficiency. Our work assumes that the database belongs to the server. The work of [31] considers approximate NNS problem in a setting very similar to ours. They focus on a biological application, which requires NNS with respect to the *edit distance*. The number of points in their dataset is relatively small (at most several thousand), so the top- $k$  selection can be done in a straightforward way (using  $O(nk)$  comparisons). We explore the different regime for the NNS problem, which is arguably even more relevant for practice: the number of points  $n$  is large (tens of thousands, millions or even more), the dimension  $d$  is not too high (several hundreds), and the distance of choice is Euclidean (for instance, by now standard and very popular NNS benchmarks [32] all fall into it). In this regime, as it turns out, the top- $k$  computation is a vast bottleneck.<sup>2</sup>

The work [33] implements the entire  $k$ -NNS computation using garbled circuits, which results in prohibitive network communication unless the dimension  $d$  is small (besides, they consider Hamming distance which is much more garbled circuit-friendly than the Euclidean distance<sup>3</sup>). The work of [34] provides a secure  $k$ -NNS solution based on the BMR protocol [35] in the *multiparty* setting where the database is

<sup>2</sup>Interestingly, when trying to implement (insecure) NNS on a GPU, the top- $k$  computation is a bottleneck as well [15].

<sup>3</sup>Since the Euclidean distance requires *multiplications*, which are known to be expensive in terms of the number of gates

distributed among different parties and another party wants to find the  $k$  nearest neighbors among all databases. Finally, the work [36] provides an extremely efficient secure NNS protocol in a different security model in which several clients use a specific hash functions and store hashes of their data on an *untrusted* server. The scheme introduces a trade-off between the search quality and an upper bound on the information leakage from hashes. In contrast, our protocols avoid any information leakage beyond the search result and the hyperparameters.

We note that the idea of combining homomorphic encryption and garbled circuits has appeared before: for instance, Gazelle [11] uses it for efficient secure neural network inference.

## C. Organization.

In Section II, we recall some background information on the cryptographic primitives used in this work. We introduce our plaintext  $k$ -NNS algorithms in Section III and the corresponding secure protocols in section IV. We present implementation details and performance results in Section V. Finally, we conclude with discussions of future directions in Section VI.

## II. PRELIMINARIES

### A. Oblivious RAM

As we have discussed, previous solutions for secure  $k$ -NNS require computing distance between the query point and all points in the database. This linear complexity is undesirable, in particular for large databases. In fact, this problem is ubiquitous in secure computation involving large datasets: existing secure computation techniques only handle computation in the circuit model, whereas in practice, many computation are efficient in the RAM model and a direct translation of RAM programs into circuits may incur large overhead.

One of the constructions that we use in this work in order to achieve sub-linear search complexity is Oblivious RAM (ORAM). ORAM was first proposed by Goldreich and Ostrovsky [37] to allow a client to outsource data storage to a remote server, and later perform efficient read/write operations without leaking access patterns. In other words, an ORAM scheme transforms the address of the block that is going to be read/written to a series of addresses that look random to the server from which server cannot discover the real address. Per each access to the database of size  $n$ , client needs to perform  $O(\text{poly}(\log n))$  accesses to the database held by the server in order to keep the real address secret. Significant research has been done to create more efficient constructions that improve the server’s memory overhead, client’s storage size, and communication between them, achieving polylogarithmic complexity in all three aspects at the same time [38], [39], [40], [41].

The idea of ORAM can also be used in the context of *secure computation*. In this scenario, the address is secret-shared among two parties that are executing the secure computation protocol and neither client nor the server know the value of the address. Here, the goal is to retrieve data at a secret

address and use it in the secure computation protocol without revealing address or data to either party. In this context, the role of ORAM client is emulated inside the secure computation protocol that produces series of physical addresses to access the database. This version of ORAM is called Distributed ORAM (DORAM), since the database is shared among two parties. In this paper, we use DORAM to achieve *sub-linear* complexity for reading from an array inside the secure computation protocol.

Wang et al. [42] suggest that a more relevant measure of the efficiency of DORAM is the *circuit complexity* of client’s functionality. The reason is that this functionality (represented as a circuit) is executed inside the secure computation protocol and is usually the most expensive part of DORAM. Therefore, many constructions have been proposed that reduce the circuit complexity at the cost of memory overhead and communication [42], [43], [44], [45]. One of the most efficient DORAM constructions, called *Floram*, is proposed by Doerner and Shelat [45]. In this work we use *Floram* in *read-only* mode, which further enhances the performance. At high-level, we implement and use two sub-routines for DORAM:

- $\text{DORAM.Init}(1^\lambda, k_A, k_B, DB) \rightarrow \overline{DB}$ . This initialization step creates an encrypted version of the database ( $\overline{DB}$ ) from the plaintext version ( $DB$ ) given two secret keys  $k_A$  and  $k_B$  (one from each party) for security parameter  $\lambda$ .
- $\text{DORAM.Read}(\overline{DB}, k_A, k_B, i_A, i_B) \rightarrow DB[i]_A, DB[i]_B$ . This subroutine performs the read operation where address  $i$  is XOR-shared between two parties as  $i_A \oplus i_B = i$ . Both parties acquire an XOR-share of the database content at address  $i$  ( $DB[i]_A$  and  $DB[i]_B$ ).

In Section IV-C, we elaborate on this construction as well as the corresponding optimizations.

### B. Additive homomorphic encryption (AHE)

In this work, we use a lattice-based additive homomorphic encryption (AHE) scheme to securely compute the Euclidean distance between two points. For our purposes, it suffices to use a private key AHE scheme, consisting of the following randomized algorithms:

- $\text{AHE.KeyGen}(1^\lambda) \rightarrow sk$ . Given security parameter  $\lambda$ , generates a secret key used for encryption and decryption.
- $\text{AHE.Enc}(sk, m) \rightarrow c$ . Encrypt a message  $m$  to a ciphertext  $c$ .
- $\text{AHE.Add}(c_1, c_2) \rightarrow c_3$ . Given encryptions of  $m_1, m_2$ , output an encryption of  $m_1 + m_2$ .
- $\text{AHE.CMult}(c, \mu) \rightarrow c'$ . Given an encryption of  $m$  and a scalar  $\mu$ , return an encryption of  $m \cdot \mu$ .
- $\text{AHE.CAdd}(c, m') \rightarrow c'$ . Given an encryption of  $m$  and a scalar  $m'$ , return an encryption of  $m + m'$ .
- $\text{AHE.Dec}(sk, c) \rightarrow m$ . Decrypt the plaintext message  $m$ .

We require our AHE scheme to satisfy standard correctness and two security properties: IND-CPA security and *function privacy*, which means that a ciphertext generated from Add and CMult operations should not reveal any information to the secret key owner, other than its underlying plaintext message.

### C. Garbled Circuits

Garbled circuit (GC) is a technique first proposed by Yao in [2] for achieving generic secure two-party computation. In this setting, two parties called *garbler* and *evaluator* hold private inputs  $x_1$  and  $x_2$  and they wish to evaluate a function  $f$  on these inputs. At the end of the evaluation, only the output  $f(x_1, x_2)$  is revealed to one or both of the parties, and no party should learn any other information about the other party’s input. Garbled circuit of a given Boolean circuit  $f$  is a triple  $(F, e, d)$ , where  $F$  is the garbled circuit,  $e$  is encoding information and  $d$  is the decoding information. Garbled circuit allows to perform secure computation as follows. The garbler chooses  $(F, e, d)$  and uses  $e$  to encode its  $x_1$  as  $X_1 = e(x_1)$ . Then, the two parties execute an *oblivious transfer* so that the evaluator obtains  $X_2 = e(x_2)$ , but the garbler learns nothing about  $x_2$ . Finally, the garbler sends the evaluator  $F, d, X_1$ . After that the evaluator evaluates encoded output  $Y = F(X_1, X_2)$  and decodes it to the final output  $y = d(Y)$  of the computation. Correctness and security of this protocol follows from correctness and security properties of garbled circuits and oblivious transfer.

Many improvements to GC have been proposed in literature, such as “free XORs” [46], meaning that XOR gates do not require the garbler to send the corresponding part of the garbled circuit, and “half-gates”, meaning that garbled AND gates are twice as small as in the original garbled circuit of Yao [47]. In addition, we use the fixed-key block cipher optimization for garbling and evaluation [48]. Using Advanced Encryption Standard (AES) as the block cipher, we leverage Intel processors’ AES instructions to perform faster garbling and evaluation.

### D. $k$ -means clustering

One of our algorithms uses the  $k$ -means clustering algorithm [49] as a subroutine. It is a simple heuristic, which finds a clustering  $X = C_1 \cup C_2 \cup \dots \cup C_k$  into disjoint subsets  $C_i \subseteq X$ , and centers  $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ , which approximately minimizes the following objective function:

$$\sum_{i=1}^k \sum_{x \in C_i} \|c_i - x\|_2^2.$$

It is immediate that for a given cluster  $C_i$ , the optimal choice of center is the *mean* of the points in  $C_i$ .

$k$ -means clustering is implemented with repeated Lloyd iterations [49] as follows. The cluster centers  $\{c_i\}$  are randomly initialized in the beginning. During each iteration, each point is attached to the nearest center with respect to the Euclidean distance. Cluster centers are recalculated at the end of an iteration by averaging over the points in each cluster. The algorithm stops either when the center-assignments converge, or when a preset maximum number of iterations is reached.

## III. PLAINTEXT $k$ -NNS ALGORITHMS

### A. High-level overview

In this paper, we present efficient and secure implementations of the following two algorithms.

a) *Algorithm 1*: The first algorithm is a heavily optimized implementation of the straightforward linear scan: we compute distances from the query point to *all* the data points, and then (approximately) select  $k_{\text{nn}}$  data points closest to the query. At a high level, we will implement distance computation using AHE, while the selection step is done using garbled circuits.

To speed up this protocol, we employ the following optimization. Computing top- $k$  naively would require a circuit consisting of  $O(nk)$  comparators. Instead, we use a special algorithm for an approximate selection of top- $k$ , which allows for a smaller circuit size (see section III-C) and will help us later when we implement the top- $k$  selection securely using garbled circuits.

b) *Algorithm 2*: The second algorithm is based on the  $k$ -means clustering (see Section II-D) and, unlike the first one, has *sublinear* query time. We now give a simplified version of the algorithm, and in Section III-B we explain why this simplified version is inadequate and provide a full description that leads to efficient implementation.

At a high level, we first compute  $k$ -means clustering of the server's dataset with  $k = k_c$  clusters. Each cluster  $1 \leq i \leq k_c$  is associated with its *center*  $c_i \in \mathbb{R}^d$ . During the query stage, we compute  $1 \leq u \leq k_c$  centers that are closest to the query, where  $u$  is a parameter to be chosen. Then we compute  $k_{\text{nn}}$  data points from the corresponding  $u$  clusters that are closest to the query, and return IDs of these points as a final answer.

Let us note that the question of choosing the hyperparameters  $k_c$  and  $u$  is fairly delicate and needs to be done separately for each particular dataset. Setting  $u$  too low leads to low accuracy, while too high values will lead to high query time. Setting  $k_c$  too low or too high leads to high query times (provided that  $u$  is tuned to achieve the desired level of accuracy).

## B. Balanced clustering and stash

To implement Algorithm 2 above (as described in Section III-A) securely without linear cost, we use secure distributed ORAM for retrieval of clusters using. In order to prevent leaking the size of each cluster, we need to set the memory block size equal to the size of the *largest* cluster in the clustering. This can be very inefficient, if the clustering at hand is not very balanced, i.e., the largest cluster is much larger than a *typical* cluster. Unfortunately, this is exactly the case in our experiments. Thus, we need a mechanism to mitigate imbalance of clusterings. Below we describe one such approach, which constitutes the *actual* version of Algorithm 2 we securely implement. With cluster balancing, our experiments achieve  $3.3\times$  to  $5.95\times$  reduction of maximum cluster sizes in different datasets.

We start with specifying the desired largest cluster size  $1 \leq m \leq n$  and an auxiliary parameter  $0 < \alpha < 1$ , where  $n$  denotes the total number of data points. Then, we find the smallest  $k$  (recall  $k$  denotes the number of centers) such that in the clustering of the dataset  $X$  found by the  $k$ -means clustering algorithm at most  $\alpha$ -fraction of the dataset lies in clusters of

size more than  $m$ . Then we consider all the points that belong to the said large clusters, which we denote by  $X'$ , setting  $n' = |X'| \leq \alpha n$ , and apply the same procedure recursively to  $X'$ . Specifically, we find the smallest  $k$  such that the  $k$ -means clustering of  $X'$  leaves at most  $\alpha n'$  points in clusters of size more than  $m$ . We then cluster these points. The algorithm terminates when every cluster has size  $\leq m$ .

At the end of the algorithm, we have  $\tilde{T}$  groups of clusters that correspond to disjoint subsets of the dataset (as a side remark, we note that one always has  $\tilde{T} \leq \log_{1/\alpha} n$ ). We denote the number of clusters in the  $i$ -th group by  $k_c^i$ , the clusters themselves by  $C_1^i, C_2^i, \dots, C_{k_c^i}^i \subseteq X$  and their centers by  $c_1^i, c_2^i, \dots, c_{k_c^i}^i \in \mathbb{R}^d$ . During the query stage, we find  $u^i$  clusters from the  $i$ -th group with the centers closest to the query point, then we retrieve all the data points from the corresponding  $\sum_{i=1}^{\tilde{T}} u^i$  clusters, and finally from these retrieved points we select  $k_{\text{nn}}$  data points that are closest to the query.

We now describe one further optimization that helps to speed up the resulting  $k$ -NNS algorithm even more. Namely, we collapse several last groups into a single set of points, which we call a *stash*. Unlike clusters from the remaining groups, we perform *linear scan* on the stash. We denote  $s$  the stash size and,  $T$  the number of remaining groups that are not collapsed.

The motivation for introducing stash is that the last few groups are usually pretty small, so in order for them to contribute to the overall accuracy meaningfully, we need to retrieve most of the clusters from them. But this means many ORAM accesses which are less efficient than the straightforward linear scan.

Note that while the simplified version of Algorithm 2 from Section III-A is well-known and very popular in practice (see, e.g., [17], [15]), our modification of the algorithm in this section, to the best of our knowledge, is new. Let us reiterate that the clustering-based  $k$ -NNS algorithms are *not* the fastest on the CPU<sup>4</sup>, but they are a perfect match for secure computation.

## C. Approximate top- $k$ selection

In both of our algorithms, we rely extensively on the following *top- $k$  selection* subroutine: given a list of  $n$  numbers  $x_1, x_2, \dots, x_n$ , find  $k \leq n$  smallest list elements in the sorted order. Let us denote the corresponding function, which outputs a tuple of size  $k$ , by  $\text{MIN}_n^k(x_1, x_2, \dots, x_n)$ . In the RAM model, computing  $\text{MIN}_n^k$  is a well-studied problem, and it is by now a standard fact that it can be computed in time  $O(n + k \log k)$  [50]. However, to perform top- $k$  selection securely, we need to implement it as a Boolean *circuit*. Suppose that all the list elements are  $b$ -bit integers. Then the desired circuit has  $bn$  inputs and  $bk$  outputs. To improve efficiency, it is desirable to design a circuit for  $\text{MIN}_n^k$  with as few gates as possible.

<sup>4</sup>However, they are known to be extremely efficient on GPU [15] due to reasons similar to the ones considered in this paper.

a) *The naïve construction:* A naïve circuit for  $\text{MIN}_n^k$  has  $O(nk)$  comparisons and hence  $O(bnk)$  gates. Roughly, it keeps a size- $k$  sorted array of the current  $k$  minima. For each value  $x_i$ , it uses a “for” loop to insert  $x_i$  into its correct location in the array, and discards the last item to keep it of size  $k$ .

b) *Sorting networks:* Another approach is to first sort the array and then take the first  $k$  elements. We could use a sorting network such as AKS [51], with  $O(bn \log n)$  gates, which is better than the naïve bound whenever  $k \gg \log n$ .

The number of gates can be further reduced to  $O(bn \log k)$  by splitting the input array into subsets of size  $k$ , and then repeatedly merging two subsets into one of size  $k$  consisting of the  $k$  smallest elements from the union of the two arrays. The merge operation can be done in  $O(bk \log k)$  gates using the AKS sorting network, and we need to perform it  $O(n/k)$  times, which gives a total of  $O(nk \log k)$  gates. This is asymptotically better than the naïve method for any super-constant value of  $k$ . However, the constant factor for AKS sorting network is prohibitively high.

On the other hand, Batchner’s sorting network [52] has slightly worse complexity  $O(n \log^2 n)$  with small constant factor. Plugging it into the above construction yields  $O(bn \log^2 k)$  gates, which is better than the naïve approach for large  $k$ .

c) *Approximate randomized selection:* We are not aware of any circuit for with  $\text{MIN}_n^k$  with  $O(bn)$  gates unless  $k$  is a constant (such bound would have been optimal, since the input has  $bn$  bits). Instead, we propose a *randomized* construction of a circuit with  $O(bn)$  gates which outputs the true top- $k$  elements with high probability. We start with shuffling the inputs in a *uniformly random order*. Namely, instead of  $x_1, x_2, \dots, x_n$ , we consider the list  $x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(n)}$ , where  $\pi$  is a uniformly (pseudo-)random permutation of  $\{1, 2, \dots, n\}$ . We require the output to be “approximately correct” (more on the precise definitions later) with high probability over  $\pi$  for every particular list  $x_1, x_2, \dots, x_n$ .

We proceed by partitioning the input list into  $l \leq n$  bins of size  $n/l$  as follows:

$$\begin{aligned} U_1 &= \{x_{\pi(1)}, \dots, x_{\pi(n/l)}\}, \\ U_2 &= \{x_{\pi(n/l+1)}, \dots, x_{\pi(2n/l)}\}, \\ &\dots, \\ U_l &= \{x_{\pi((l-1)n/l+1)}, \dots, x_{\pi(n)}\}. \end{aligned}$$

Our circuit works in two stages: first we compute a minimum within each bin  $M_i = \min_{x \in U_i} x$ , and then we output  $\text{MIN}_l^k(M_1, M_2, \dots, M_l)$  as a final result using the naïve circuit for  $\text{MIN}_l^k$ . The circuit size is  $O(b \cdot (n + kl))$ , which is  $O(bn)$  if  $kl = O(n)$ .

Intuitively, if we set  $l$  to be large enough, the above circuit outputs a “good” answer with high probability over  $\pi$ . We state and prove two theorems formalizing the notion of “good” in two different ways.

**Theorem 1.** *For every  $n$ ,  $0 < \delta < 1/2$  and  $k \ll_\delta \sqrt{n}$ , there exists  $l = O(k/\delta)$  such that the intersection of the output of*

*our circuit with  $\text{MIN}_n^k(x_1, x_2, \dots, x_n)$  is at least  $(1 - \delta)k$  entries in expectation over the choice of  $\pi$ .*

This yields a circuit of size  $O(b \cdot (n + k^2/\delta))$ .

*Proof.* By the main result of [53] and the bound  $k \ll_\delta \sqrt{n}$ , it is enough to prove the bound on the expected size of the required intersection for the case when we sample a bin for each element *independently*. For this sampling model, it is not hard to see that the desired expectation is

$$l \cdot \left(1 - \left(1 - \frac{1}{l}\right)^k\right),$$

which is at least  $(1 - \delta)k$  for  $l \sim k/\delta$ .  $\square$

**Theorem 2.** *For every  $n$  and  $0 < \delta < 1/2$  and  $k \ll_\delta n^{1/3}$ , there exists  $l = O(k^2/\delta)$  such that the output of the circuit is exactly  $\text{MIN}_n^k(x_1, x_2, \dots, x_n)$  with probability at least  $1 - \delta$  over the choice of  $\pi$ .*

This yields a circuit of size  $O(b \cdot (n + k^3/\delta))$ , which is worse than the previous bound, but the corresponding correctness guarantee is stronger.

*Proof.* Similar to the proof of the previous theorem, it is enough to show the result for the independent sampling of bins. But for such sampling, the result holds by the birthday paradox.  $\square$

In some applications, it is enough to output a binary vector of length  $n$  with exactly  $k$  ones on the positions that correspond to the  $k$  smallest entries of the list. It was known how to do this in  $O(b^2n)$  gates [54], and we show how to improve this to the optimal  $O(bn)$  gates. Such a circuit can be used for the linear scan, but for the clustering-based algorithm, we need to return  $k$  smallest entries explicitly. Due to this requirement and also the fact that the new  $O(bn)$ -sized circuit has a higher hidden constant than the above randomized construction, we decided not to implement it. For completeness and for the future reference, we describe the new circuit in Section D.

#### D. Approximate distances

To speed up the top- $k$  selection further, instead of exact distances, we will be using *approximate* distances. Namely, instead of storing full  $b$ -bit distances, we discard  $r$  low-order bits, and the overall number of gates in the selection circuit becomes  $O((b - r) \cdot (n + kl))$ .

For the clustering-based algorithm, we set  $r$  differently when we select closest cluster centers and when we select closest data points, which allows for a more fine-grained parameter tuning.

#### E. Putting it together

We now give a high-level summary of our algorithms. For the linear scan, we use the approximate top- $k$  selection to return the  $k_{\text{nn}}$  IDs after computing distance between query and all points in the database; for the clustering-based algorithm, we use approximate top- $k$  for retrieving  $u^i$  clusters in  $i$ -th group for all  $i \in \{1, \dots, T\}$ . Then, we compute the closest

$k_{\text{nn}}$  points from query to all the retrieved points. Meanwhile, we compute the approximate top- $k$  with  $k = k_{\text{nn}}$  between query and the stash. Finally, we compute and output the  $k_{\text{nn}}$  closest points from the above  $2k_{\text{nn}}$  candidate points.

Note that in the clustering-based algorithm, we use exact top- $k$  selection for retrieved points and approximate selection for cluster centers and stash. The main reason is that the approximate selection requires a random shuffle of the input points. This shuffle must be known only to the server and not to the client to ensure that there is no additional information leakage when the algorithm is implemented securely. Jumping ahead to the secure protocol in the next section, the points we retrieve from the clusters will be in secret-shared format. Thus, doing approximate selection on retrieved points would require a secure two-party shuffling protocol, which is expensive. Therefore, we run a naïve circuit for exact computation of top- $k$  for the retrieved points.

#### F. What does the output of our algorithms leak?

We briefly discuss the potential leakage from the *output* of our approximate  $k$ -NNS algorithms as described in Section III. Note that this discussion is independent from the privacy guarantee of our secure protocols, which have no leakage themselves. Here, we only investigate what is leaked to the client via the output.

First, we discuss the leakage of our algorithms beyond the exact  $k$ -NNS functionality. Note that in our algorithms the final  $k$  points sent back to the client are not fixed – they will be affected by the random shuffles done by the server. The client can exploit this by asking the same query many times. In this case, the client can count the number of times each point is returned. In our linear scan algorithm, the closer a point is to the query, the more likely it is to be returned. Therefore, from the point counts, the client could reconstruct the top- $k$  points (that is, the result of exact  $k$ -NNS) with high probability. Moreover, the client could also get hints about the top- $(k+1)$ ,  $k+2$ ,  $\dots$  points. However, we conjecture that the leakage is small in practice. In fact, for a slightly different way of approximate selection, where instead of putting  $n$  points into  $l$  equal-sized buckets, we assign each point uniformly to one of the  $l$  buckets, we verified the following result: suppose a point has distance rank  $k+r$  from the query, then the probability that the point is included in the result of our approximate  $k$ -NNS algorithms decays exponentially with  $r$ . In other words, points which are far away from the query have very little chance of appearing from the list.

Also, there is a line of work analyzing the leakage of exact  $k$ -NNS. For example, [55] shows that in low-dimensional databases, one can approximately reconstruct the database after issuing sufficiently many  $k$ -NNS queries. Here, we remark that the current techniques of database recovery from  $k$ -NNS query results still do not scale well to large-dimension data considered in this work.

We also note that by asking many queries adaptively, the client can recover an approximate  $k$ -NN graph of the dataset, which contains lots of valuable information about the data,

including community structure. To prevent this, one needs to restrict the client in the number of query and the degree of adaptivity.

## IV. SECURE PROTOCOL FOR $k$ -NNS

### A. Overview of our protocol

We give a high-level overview of our secure protocols implementing the functionalities in the previous section, followed by description of individual subroutines (AHE, garbled circuits, and DORAM). We start with the clustering-based protocol. First, we give an illustration of the protocol in Figure 1.

Recall that the server’s input to the protocol is a partition of its database  $X$  into clusters and a stash

$$X = \left[ \bigcup_{i=1}^T \bigcup_{j=1}^{k_c^i} C_j^i \right] \cup S,$$

such that each cluster  $C_j^i$  has size at most  $m$ . Let  $c_j^i$  denote the center of  $C_j^i$ . Our protocol works in the following stages:

**Setup.** The server and the client execute DORAM.Init to insert all clusters from all groups into DORAM, with one cluster in each block. Clusters are padded by “infinitely far points” if necessary to reach size  $m$ .

**Query.** This stage consists of the following steps.

- 1) The server performs two independent random shuffles on the cluster centers and stash points necessary for the approximate top- $k$ .
- 2) For each  $i \in \{1, \dots, T\}$ ,
  - The client and server use AHE to compute secret shares of  $d_j^i = \|\mathbf{q} - c_j^i\|_2^2$  for all  $j$ .
  - Client and server run approximate top- $k$  selection algorithm from Section III-C using garbled circuits, with  $k = u_i$ , and output secret shares of  $u_i$  cluster indices.
- 3) Client and server input the secret shares of the  $u_{\text{all}} = \sum_{i=1}^T u_i$  indices  $(i_1, j_1), \dots, (i_{u_{\text{all}}}, j_{u_{\text{all}}})$  obtained in previous step into DORAM.Read to retrieve all points in  $C := C_{j_1}^{i_1} \cup \dots \cup C_{j_{u_{\text{all}}}}^{i_{u_{\text{all}}}}$  in secret shared form.
- 4) Use AHE to compute secret shares of distances between  $\mathbf{q}$  and all points in  $C \cup S$ .
- 5) Use garbled circuit to securely evaluate a naïve top- $k$  circuit, compute secret shares of IDs and distances of  $k_{\text{nn}}$  closest points in  $C$  to query.
- 6) Use garbled circuit to securely evaluate the approximate top- $k$  circuit from Section III-C to compute secret shares of IDs and distances of  $k_{\text{nn}}$  closest points in  $S$  to query.
- 7) Use garbled circuit to evaluate the naïve top- $k$  selection circuit which take as input secret shares of the above  $2k_{\text{nn}}$  points (and distances) and output the IDs of closets  $k_{\text{nn}}$  points to the client.

Now, our linear scan protocol can be obtained by setting the stash equal to the entire database, i.e.  $S = X$ , and skipping the clustering and DORAM altogether, so the setup phase only requires server to randomly shuffle all points. Then, we

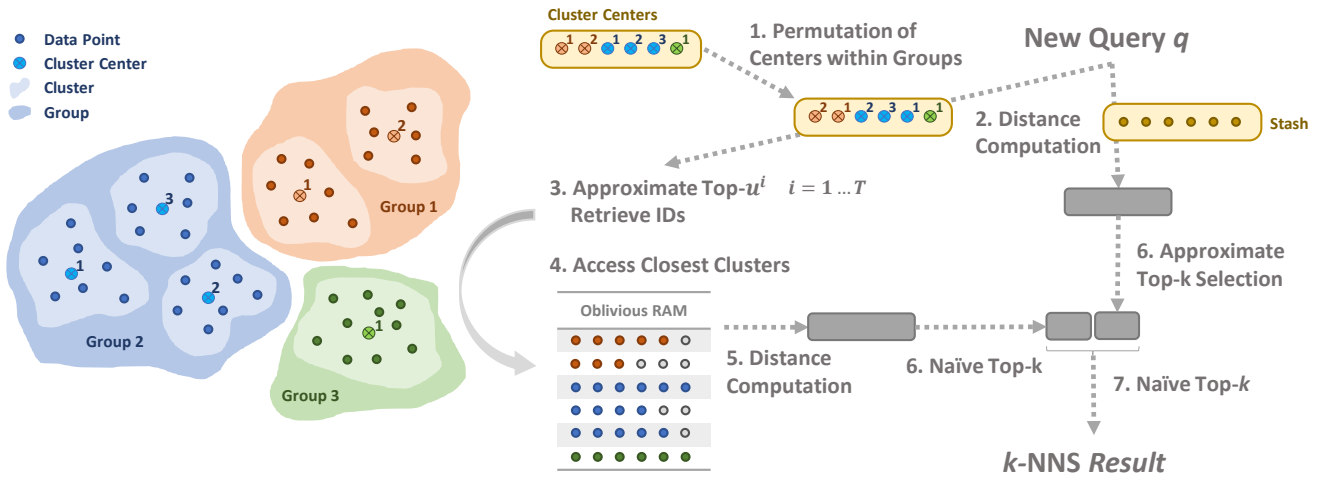


Fig. 1: Global computation and data flow of SANNS.

execute step (3) and (5) in the above query procedure to obtain a list of  $k_{mn}$  IDs.

### B. Distance computation from AHE

It is well-known that secure computation of Euclidean distances can be done using AHE. Among the existing AHE schemes, we select the lattice-based Brakerski/Fan-Vercauteren (BFV) scheme [56], [7] with the nice property that it supports efficient single-instruction-multiple-data (SIMD) operations on encrypted vectors. This allows us to compute distances from the query point to many points of the dataset at once. The idea of using the BFV scheme to perform fast secure linear operations is in the same spirit as [11]. However, compared to [11], our approach avoids expensive ciphertext *rotations*. Also, we modify the SIMD encoding technique to fit our scenario, notably removing the restriction on the plaintext modulus and perform computation modulo a power of two instead. The benefit of computation modulo powers of two (as opposed to modulo prime  $p$ ) is that it allows us to later avoid a costly addition modulo  $p$  transformation inside a garbled circuit when reconstructing distances from secret shares. Thus, our approach is more efficient and more compatible with the garbled circuit components of our protocols.

More precisely, in order to enable SIMD operation such as elementwise multiplication of vectors in the BFV scheme, we need to work with plaintexts consisting of integers modulo some prime  $p \equiv 1 \pmod{2N}$ , where  $N$  is the ring dimension parameter. However, we observe that our distance computation protocol only requires efficient multiplication between *scalars* and vectors. Therefore we can drop the requirement on the plaintext modulus and perform computation modulo some power of two without losing efficiency.

We now describe our distance computation protocol in more detail. Recall that plaintext space of the BFV scheme is a polynomial ring  $R_t := \mathbb{Z}_t[x]/(x^N + 1)$ , where we take  $N$  to be a power of 2 and  $t$  an integer modulus. So a plaintext is represented as a polynomial with degree less than  $N$  and coefficients in  $\mathbb{Z}_t$ . Then, the client encodes each coordinate of

the query separately into the constant coefficient in  $R_t$ : if the query is  $\mathbf{q} = (q_1, \dots, q_d) \in \mathbb{R}^d$ , then we encode it as

$$f_i = q_i + 0 \cdot x + \dots + 0 \cdot x^{N-1}, \quad 1 \leq i \leq d.$$

For the sake of simplicity, assume that the server has  $N$  points  $\mathbf{p}_1, \dots, \mathbf{p}_N$ . Then, it will use  $d$  plaintexts, each encoding one coordinate of all points, resulting in

$$g_i = p_{1,i} + p_{2,i}x + \dots + p_{N,i}x^{N-1}, \quad 1 \leq i \leq d.$$

Then, we could verify that

$$\sum_{i=1}^d f_i g_i = \sum_{j=1}^N \langle \mathbf{q}, \mathbf{p}_j \rangle x^{j-1}.$$

That is, we could compute  $N$  dot products using  $d$  homomorphic scalar multiplications and additions. Our protocol works by letting the client encrypt each  $f_i$  into a ciphertext  $c_i$  and send to the server; then the server uses  $\text{AHE.CMult}$  and  $\text{AHE.Add}$  to compute a ciphertext encrypting  $h(x) = \sum_{j=1}^N \langle \mathbf{q}, \mathbf{p}_j \rangle x^{j-1}$ . The server then samples a random polynomial  $r(x)$  and uses  $\text{AHE.CAdd}$  to compute encryption of  $h(x) + r(x)$ , which it sends back to the client. The client then decrypts the ciphertext to obtain  $h(x) + r(x)$ , and the server keeps  $r(x)$ ; in other words, the client and the server hold secret shares of  $\langle \mathbf{q}, \mathbf{p}_j \rangle$  modulo  $t$ . Then, secret shares of Euclidean distances can be reconstructed via local operations, using the following identity

$$\|\mathbf{q} - \mathbf{p}_j\|_2^2 = \|\mathbf{q}\|_2^2 - 2\langle \mathbf{q}, \mathbf{p}_j \rangle + \|\mathbf{p}_j\|_2^2.$$

We need to slightly modify the above routine when computing distances of points retrieved from DORAM. Here the server does not hold points in the clear: instead, the client and server secret share the points and their squared Euclidean norms. We use  $\langle x \rangle_C$  and  $\langle x \rangle_S$  to denote the client and server's shares of a private input  $x$ , such that  $x = \langle x \rangle_C + \langle x \rangle_S \pmod{t}$ . Then, we only need to securely compute dot products between  $\mathbf{q}$  and each  $\langle \mathbf{p}_j \rangle_S$ , since the following holds:

$$\|\mathbf{q} - \mathbf{p}_j\|_2^2 = \|\mathbf{q}\|_2^2 - 2\langle \mathbf{q}, \langle \mathbf{p}_j \rangle_C \rangle + \|\langle \mathbf{p}_j \rangle_C\|_2^2 - 2\langle \mathbf{q}, \langle \mathbf{p}_j \rangle_S \rangle + \|\langle \mathbf{p}_j \rangle_S\|_2^2.$$



### C. Point Retrievals Using DORAM

In our work, we use *Floram*, a DORAM construction proposed by Doerner and Shelat [20] in order to securely retrieve data points in the database. Floram is one of the most efficient DORAM schemes that has a low circuit complexity at the cost of linear *local* computation by the two parties holding the secret-shared of database. While there exists DORAM schemes with logarithmic local computation, circuit complexity, and communication [43], Floram is more efficient in practice due to light-weight local computations. Here, we briefly explain the functionality of Floram and refer the reader to the original paper [20] for a more detailed description.

In Floram, both parties hold *identical* copies of the masked database. Let us denote the plaintext database as  $DB$ , word at address  $i$  as  $DB[i]$ , and the masked database as  $\overline{DB}$ . At the initialization time, parties hold a secret-shared of a database. We denote party  $A$  and  $B$ 's shares as  $DB_A$  and  $DB_B$ , respectively where  $DB[i] = DB_A[i] \oplus DB_B[i]$ . Prior to retrieving a word, parties  $A$  and  $B$  sample a random key for a Pseudo-Random Function (PRF), denoted by  $k_A$  and  $k_B$ , respectively. They bitwise-XORs each word at address  $i$  with  $PRF_{k_A}(i)$  or  $PRF_{k_B}(i)$ , send the result to the other party, and create the masked database as

$$\begin{aligned}\overline{DB}[i] &= DB_A[i] \oplus PRF_{k_A}(i) \oplus DB_B[i] \oplus PRF_{k_B}(i) \\ &= DB[i] \oplus PRF_{k_A}(i) \oplus PRF_{k_B}(i)\end{aligned}$$

Note that the bit-length of the output of  $PRF$  is equal to (or bigger than) the bit-length of each word. At a high-level, Floram's data retrieval functionality consists of the two main parts: token generation using Functional Secret Sharing (FSS) and data unmasking using PRF. FSS enables two parties to secret-share a point function. A point function  $f_{\alpha, \beta}(x)$  takes the value of  $\beta$  where  $x = \alpha$  and zero otherwise. Gilboa and Ishai [57] have shown how one can secret-share a point function with shares sub-linear with respect to function's domain. Built upon this, Boyle et al. [58], [59] constructed a two-server Private Information Retrieval (PIR) system. In Floram, FSS is used to securely generate two bit-arrays (one for each party). The input is a secret-shared target index  $\alpha$ , and output is two random-looking bit-array to each party, subject to the constraint that the XOR of which is equal to  $(f_{\alpha, 1}(i))_{i=0}^{n-1}$ . Then, by taking the dot product between the bit-array and the corresponding memory words, the parties now have a correct XOR-share of the masked word at the target index  $\alpha$ . Finally, parties unmask the word by inputting their shares as well as their PRF keys to garbled circuit (GC), where the two PRF masks are removed and the desired word is re-shared and one share is output to each party.

We have implemented Floram with the following optimizations as we describe below. The first two are proposed by the Floram paper and we integrate those in our implementation too. The third optimization reduces the overhead of FSS evaluation. We also propose to replace AES with Kreyvium to realize the functionality of the PRF. Last but not least, we reduce the number of interactions between two parties when

accessing the database at several different indices. In what follows, we discuss these optimizations in more detail.

**Constant PRG (CPRG).** The costliest part of FSS is many evaluations of Pseudo-Random Generator (PRG) in the GC protocol. Doerner and Shelat [20] propose an optimization that replaces secure evaluation of PRG with  $\log_2 n$  simple secure computations. As a result of this technique, the round complexity is increased to  $\log_2 n$  per access but all PRG evaluations (required in FSS) are performed in plaintext.

**Tree trimming.** The second optimization proposed for *read* operations in Floram is to avoid evaluating certain number of last layers in FSS tree at the cost of small computation overhead. However, the overhead is quickly paid off due to the exponential growth of the last layers in FSS tree. We refer the reader to the Floram paper [20] for more detailed explanation.

**Precomputing OT.** Recall that with CPRG technique, two parties have to execute the GC protocol  $\log_2 n$  times iteratively which in turn requires  $\log_2 n$  set of Oblivious Transfers (OTs). Performing consecutive OTs can significantly slow down the FSS evaluation. In order to mitigate the overhead, we propose to use Beaver OT precomputation protocol [60] which enables us to perform all necessary OTs on random values in the beginning of FSS evaluation with a very small communication for each GC invocation.

**Kreyvium as PRF.** In original Floram, PRF is implemented using Advanced Encryption Standard (AES). While computing AES is fast in plaintext due to Intel AES new instructions, it requires many AND gates to be evaluated within a garbled circuit. Thus, we propose a more efficient solution based on stream ciphers. In particular, we implement our PRF using Kreyvium [21] which requires significantly fewer number of AND gates (see Appendix C for various related trade-offs). However, evaluating Kreyvium in plaintext during the initial database masking adds more overhead compared to AES. To mitigate the overhead, we pack multiple (512 in our case) invocations of Kreyvium and evaluate them simultaneously by using Advanced Vector Extensions (AVX-512) instructions provided by Intel processors.

**Multi-address access.** All of the aforementioned optimizations improve the performance of a single access. Accessing the database at  $k$  different locations, requires  $k \log_2 n$  number of interactions. If these memory accesses are non-adaptive, then the same process can be implemented much more efficiently by fusing all of the access procedures reducing the number of rounds to merely  $\log_2 n$ .

### D. Top- $k$ selection using Garbled Circuits

We start with secret shares of distances modulo  $2^{b_d}$ , where  $b_d$  is the number of bits required to store a single distance. At a high level, we implement the top- $k$  selection by plugging in the randomized circuit described in Section III-C into Yao's garbled circuits framework [2]. There are some further optimizations we make in order to improve the performance.

First, instead of working with exact distances, we round them, which allows us to reduce the circuit size significantly

(see Section III-D). This is done by simply discarding some lower order bits after adding the secret shares in the garbled circuit.

The second optimization comes from the implementation side. Using generic MPC frameworks such as ABY [28] ends up being problematic for us, since such frameworks require to store the whole circuit explicitly with accompanying bloated data structures. However, our top- $k$  circuit is highly structured (i.e., it is a composition of a certain small circuit with itself many times), which allows us to work with it “locally”. This means that the memory consumption of the garbling and the evaluation algorithms is essentially independent  $n$ , which makes them much more cache-efficient, and, as a result, much faster.

For this, we use our own GC implementation with most of the standard optimizations [35], [61], [48], [62]<sup>5</sup>, which allows us to save more than an order of magnitude in both time and memory usage compared to ABY.

## V. IMPLEMENTATION AND PERFORMANCE RESULTS

### A. Environment

We perform the evaluation on two Azure F72s\_v2 instances (with 72 virtual CPUs and 144 Gb of RAM each) hosted in the “West US 2” availability zone. We evaluate our algorithms in 1 and 72 threads (for the query procedure, the preprocessing and OT phases are always single-thread). We implement networking using ZeroMQ: the latency between instances ends up being around 0.5 ms, while the throughput ranges between 374 MB/s on a single thread and 3.30 GB/s on 72 threads. We also perform an experiment on two instances hosted in “West US 2” and “East US” availability zones. In that case, the networking is a good deal slower: the latency is 34 ms and the throughput ranges between 36 MB/s for a single thread, and 2.0 GB/s for 72 threads. We use g++ 7.3.0, Ubuntu 18.04, SEAL 2.3.1 [19] and libOTe [63] for the OT phase (in the single-thread mode). We implement balanced clustering as described in Section III-B using PyTorch and run it on four NVIDIA Tesla V100 GPUs. It is done once per dataset and takes several hours (with the bottleneck being the vanilla  $k$ -means clustering described in Section II-D).

### B. Datasets

We evaluate our algorithms as well as baselines on two datasets: SIFT ( $n = 1\,000\,000$ ,  $d = 128$ ) is a standard dataset of image descriptors [22] that can be used to compute similarity between images; Deep1B ( $n = 1\,000\,000\,000$ ,  $d = 96$ ) is also a dataset of image descriptors [23], which is more modern and are feature vectors obtained by passing images through a deep neural network (for more details see the original paper [23]). We conduct the evaluation on two subsets of Deep1B that consist of the first 1 000 000 and 10 000 000 images, which we label Deep1B-1M and Deep1B-10M, respectively. SIFT comes with 10 000 sample queries

which we use for evaluation; for Deep1B-1M and Deep1B-10M, we use a sample of 10 000 data points, which we remove from the dataset, as queries. For all the datasets we use Euclidean distance to measure similarity between points. Note that the Deep1B-1M and Deep1B-10M datasets are normalized to lie on the unit sphere.

Note that all of the above datasets have been extensively used in nearest neighbors benchmarks. In particular, SIFT is a part of ANN Benchmarks [64], where a large array of NNS algorithms has been thoroughly evaluated. Deep1B has been used for evaluation of NNS algorithms in [23], [15], [16] and a number of other papers.

### C. Accuracy

In our experiments, we require the algorithms to return 10 nearest neighbors and measure accuracy as the average of the number of correctly returned points over the set of queries (we refer to this later as “10-NN accuracy”). We evaluate our algorithms requiring that the 10-NN accuracy is at least 0.9, which is a level of accuracy considered to be acceptable in practice.

### D. Quantization of coordinates

For SIFT, coordinates of points and queries are already small integers between 0 and 255, so we leave them as is. For Deep1B, the coordinates are real numbers, and we quantize them to 8-bit integers uniformly between the minimum and the maximum coordinates for the dataset. In experiments, such quantizations barely affect the 10-NN accuracy compared to using the true floating point coordinates.

### E. Cluster size balancing

As noted in Section III-B, our cluster balancing algorithm achieves the crucial bound over the maximum cluster size needed for efficient ORAM retrieval of candidate points. In our experiments, for SIFT, Deep1B-10M, and Deep1B-1M, the balancing algorithm reduced the maximum cluster size by factors of  $4.95\times$ ,  $3.67\times$ , and  $3.31\times$ , respectively.

### F. Notation

Here we list the hyperparameters used by our algorithms. See Figure 5 and Figure 6 for the values that we use for various datasets.

Main hyperparameters:

- $n$  is the number of data points
- $d$  is the dimension
- $k_{nn}$  is the number of data points we need to return as an answer
- $T$  is the number of groups
- $k_c^i$  is the total number of clusters for the  $i$ -th group,  $1 \leq i \leq T$
- $m$  is the largest cluster size
- $u^i$  is the number of closest clusters we retrieve for the  $i$ -th group,  $1 \leq i \leq T$
- $u_{all}$  is the total number of clusters we retrieve,  $u_{all} = \sum_{i=1}^T u^i$

<sup>5</sup>For oblivious transfer, we use libOTe [63]

	Threads	Algorithm	Overall query	ORAM	Top- $k$	Distances	OT phase	Preprocessing
SIFT	1	Linear scan	35.4 s 4.52 GB	None	15.6 s 4.42 GB	19.8 s 98.8 MB	2.99 s 950 MB	None
		Clustering	8.63 s 1.77 GB	4.38 s 1.07 GB	1.98 s 645 MB	2.22 s 56.7 MB	0.63 s 166 MB	12.9 s 484 MB
	72	Linear scan	6.15 s 4.52 GB	None	2.54 s 4.42 GB	3.08 s 98.8 MB	N/A	None
		Clustering	<b>2.36 s</b> 1.79 GB	0.92 s 1.07 GB	1.00 s 666 MB	0.35 s 56.7 MB	N/A	N/A
Deep1B-1M	1	Linear scan	30.0 s 4.50 GB	None	15.1 s 4.42 GB	14.9 s 86.2 MB	3.07 s 950 MB	None
		Clustering	7.44 s 1.59 GB	3.87 s 921 MB	1.86 s 621 MB	1.67 s 44.1 MB	0.59 s 153 MB	11.0 s 407 MB
	72	Linear scan	6.02 s 4.50 GB	None	2.66 s 4.42 GB	2.87 s 86.2 MB	N/A	None
		Clustering	<b>2.33 s</b> 1.61 GB	0.91 s 921 MB	1.07 s 640 MB	0.33 s 44.1 MB	N/A	N/A
Deep1B-10M	1	Linear scan	390 s 47.9 GB	None	187 s 47.4 GB	203 s 518 MB	32.6 s 10.4 GB	None
		Clustering	31.6 s 5.53 GB	18.0 s 3.12 GB	7.23 s 2.35 GB	6.33 s 59.4 MB	1.83 s 576 MB	86.3 s 3.72 GB
	72	Linear scan	75.9 s 47.9 GB	None	54.0 s 47.4 GB	17.0 s 518 MB	N/A	None
		Clustering	<b>6.37 s</b> 5.59 GB	2.94 s 3.12 GB	2.74 s 2.41 GB	0.68 s 59.4 MB	N/A	N/A

Fig. 2: Performance of our algorithms on two “West US 2” Azure instances. We show the break down of the running time and communication between the parts of the algorithm. “Overall query time” does not include the OT phase, which is done once per query. Preprocessing is done once per client. We run OT and preprocessing in a single thread. Also we measure overall query time as the maximum between server and client, but measure the parts on the server.

	Threads	Algorithm	Overall query	ORAM	Top- $k$	Distances	OT phase	Preprocessing
SIFT	1	Linear scan	130 s	None	103.7 s	24.9 s	30.2 s	None
		Clustering	61.8 s	41.3 s	16.09 s	3.56 s	4.95 s	23.6 s
	72	Linear scan	21.5 s	None	4.45 s	13.5 s	N/A	None
		Clustering	11.5 s	3.70 s	5.30 s	2.01 s	N/A	N/A
Deep1B-1M	1	Linear scan	125 s	None	104 s	20.1 s	23.9 s	None
		Clustering	47.1 s	27.6 s	16.4 s	3.09 s	4.56 s	20.2 s
	72	Linear scan	20.5 s	None	4.43 s	12.9 s	N/A	None
		Clustering	11.2 s	3.78 s	5.29 s	1.90 s	N/A	N/A
Deep1B-10M	1	Linear scan	1400 s	None	1190 s	204 s	250 s	None
		Clustering	172 s	103 s	58.3 s	10.1 s	14.5 s	165 s
	72	Linear scan	211 s	None	186 s	16.8 s	N/A	None
		Clustering	29.7 s	9.00 s	16.4 s	4.04 s	N/A	N/A

Fig. 3: Similar to Figure 2, but now the instances are hosted on “West US 2” and “East US”. so the running times are higher due to the slower networking. We do not report communication, since it’s the same to Figure 2.

- $s$  is the *stash size*
- $l^i$  is the number of bins we use to speed up the selection of closest clusters for the  $i$ -th group,  $1 \leq i \leq T$
- $l_s$  is the number of bins we use to speed up the selection of closest points for the stash
- $b_c$  is the number of bits necessary to encode one *coordinate*
- $b_d$  is the number of bits necessary to encode one *distance* ( $b_d = 2b_c + \lceil \log_2 d \rceil$ )
- $b_{cid}$  is the number of bits necessary to encode the ID of a *cluster* ( $b_{cid} = \lceil \log_2 \left( \sum_{i=1}^T k_c^i \right) \rceil$ )
- $b_{pid}$  is the number of bits necessary to encode the ID of a *point* ( $b_{pid} = \lceil \log_2 n \rceil$ )
- $r_c$  is the number of bits we discard when computing distances to *centers of clusters*,  $0 \leq r_c \leq b_d$
- $r_p$  is the number of bits we discard when computing distances to *points*,  $0 \leq r_p \leq b_d$

Additional hyperparameters:

- $\alpha$  is the allowed fraction of points in large clusters during the preprocessing
- $N$  is the ring dimension in BFV scheme;  $q$  is the ciphertext modulus.
- $t$  is the plaintext modulus in BFV scheme and the modulus for the intermediate secret sharings. Note that  $t = 2^{b_d}$ .

### G. Parameter choices

We initialized the BFV scheme with parameters  $N = 8192$ ,  $t = 2^{23}$  and a 180-bit modulus  $q$ . For the parameters such as standard deviation error and secret key distribution we use SEAL default values. These parameters allow us to use the noise flooding technique to provide 108 bits of statistical circuit privacy<sup>6</sup>. We used the LWE estimator<sup>7</sup> by Albrecht et al. [65] to estimate the security level of the scheme, which suggests 141 bits of security.

Let us describe how we set the hyperparameters of our algorithms. Both of our algorithms (especially the clustering-based) have quite a few moving parts that nontrivially affect the overall performance. See Section V-F for the full list of hyperparameters, below we list the one that affect the performance for both of our algorithms:

- Both algorithms depend on  $n$ ,  $d$ ,  $k_{nn}$ , which depend on the dataset and our requirements;
- Besides that, linear scan depends on  $l_s$ ,  $b_c$  and  $r_p$ ,
- And the clustering-based algorithm depends on  $T$ ,  $k_c^i$ ,  $m$ ,  $u^i$ ,  $s$ ,  $l^i$ ,  $l_s$ ,  $b_c$ ,  $r_c$  and  $r_p$ , where  $1 \leq i \leq T$ .

For both of the algorithms, we use the *total number of AND gates* in the top- $k$  and the ORAM circuits as a proxy for both communication and running time. Moreover, for simplicity we neglect the FSS part of ORAM, since it does not affect the performance much. We refer the reader to Section A for the exact formulas used in our cost model.

<sup>6</sup>We refer the reader to [11] for details on the noise flooding technique

<sup>7</sup>We used the most recent commit (3019847) from <https://bitbucket.org/malb/lwe-estimator>

Overall, we search for the hyperparameters that yield 10-NN accuracy at least 0.9 (approximately) minimizing the total number of AND-gates. We list in Figure 5 and Figure 6 the settings we use.

### H. Evaluation

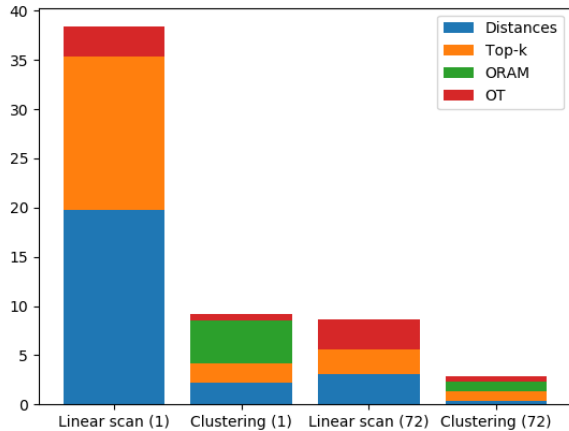
Figure 2 shows the running times and communication volumes for both of our algorithms evaluated on SIFT, Deep1B-1M and Deep1B-10M run on two “West US 2” instances. Since the OT phase and per-client preprocessing are implemented only in a single-thread regime, we mark the respective entries in the “multi-thread” rows with “N/A”. Let us note however that both of these phases should be easily parallelizable. We find that the clustering-based algorithm consistently outperforms the linear scan, both in terms of the running time and the communication, both for 1 and 72 threads. On Deep1B-10M, the gap for the respective characteristics exceeds an order of magnitude. It is interesting that for a single thread and a single query, the clustering-based algorithm beats the linear scan *even taking the per-client preprocessing time into account*. We do not report the timing of hyperparameter tuning and clustering, since it needs to be done only once per dataset. We also evaluate our algorithms on a slower network connection: between a “West US 2” and an “East US” instance, see Figure 3. The results are qualitatively similar to Figure 2.

Let us now compare the numbers we obtain with two baselines. First, we use the arithmetic mode of ABY [28] to compute distances from a query to the data points. We find that on SIFT it takes 620 seconds and 167 GB of communication, which is dramatically worse than what can be done by AHE (2.22 s, 56.7 MB). On Deep-1B-1M, ABY takes about the same time, and on Deep-1B-10M, it consumes more than all the available RAM on our instances, but it likely to be an order of magnitude slower. Second, we evaluate the naïve top- $k$  circuit that consists of  $O(nk)$  comparisons that has been used in the prior work (e.g., in [33]) using our GC implementation. On SIFT it takes 147 seconds and 24.7 GB of communication, while our better circuit takes merely 15.6 seconds and 4.42 GB of communication, improving by almost an order of magnitude. We note that the gap in communication is around 5x due to the fact that we compute distances from secret shares, which becomes one of the bottlenecks for our faster top- $k$  selection.

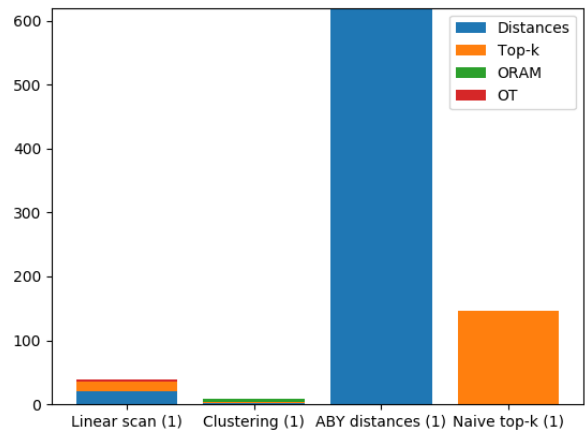
We summarize the running times of our algorithms as well as the baselines on Figure 4.

## VI. CONCLUSIONS AND FUTURE DIRECTIONS

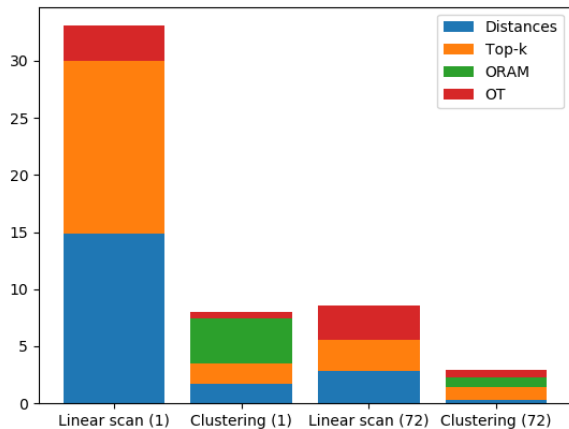
In this work, we design new secure computation protocols achieving approximate  $k$ -Nearest Neighbors Search functionality between a client with query and a server holding a database, with the Euclidean distance metric. Our solution combines several state-of-the-art cryptographic primitives such as lattice-based additively homomorphic encryption, FSS-based distributed ORAM and garbled circuits with by now standard optimizations. We also design tailored plaintext approximate  $k$ -NNS algorithms to produce good accuracy while



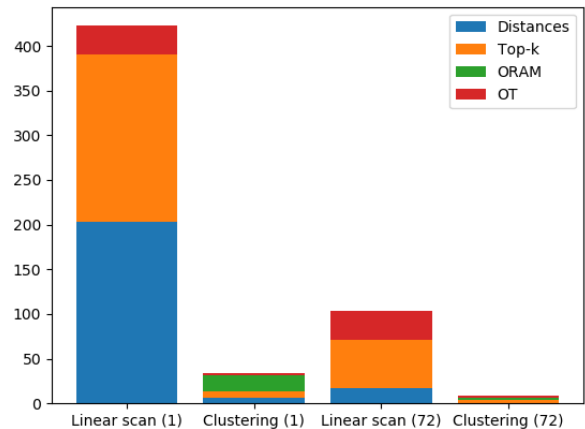
(a) Performance of our algorithms on SIFT



(b) Performance of our algorithms on SIFT next to the baselines



(c) Performance of our algorithms on Deep1B-1M



(d) Performance of our algorithms on Deep1B-10M

Fig. 4: Comparison of our algorithms run on 1 and 72 threads on two “West US 2” instances. The y-axis is the running time (in seconds). OT phase is always run single-threaded. For SIFT we compare our algorithms with distances computed in ABY as well as the naïve top- $k$  circuit.

at the same time achieve good efficiency when implemented securely. Notably, our clustering-based protocol is the first *sublinear* secure protocol for approximate  $k$ -NNS. Our performance results show that our solution scales well to massive datasets consisting of ten million points.

We highlight some directions for future works:

- Our construction can be proved secure in the honest-but-curious model, but it would be interesting to extend our protocols to protect against malicious scenarios, where the client or the server can deviate from the protocol in order to learn about the other party’s data or manipulate the output of the other party.
- Another open problem regards the PRF function. We decided to use Kreyvium instead of AES in order to

reduce communication between the parties, but when the cipher needs to be evaluated in the clear, AES is still more efficient thanks to optimized hardware implementation. It would be interesting to investigate on improvements of the plaintext implementation of Kreyvium. This could also be useful to improve its ciphertext evaluation.

- It would be interesting to implement other sublinear  $k$ -NNS algorithms securely, most notably Locality-Sensitive Hashing (LSH) [18], which has *provable* sublinear query time.

## REFERENCES

- [1] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in

- Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 586–595.
- [2] A. C.-C. Yao, “How to generate and exchange secrets,” in *Foundations of Computer Science, 1986., 27th Annual Symposium on*. IEEE, 1986, pp. 162–167.
  - [3] O. Goldreich, S. Micali, and A. Wigderson, “How to play any mental game or A completeness theorem for protocols with honest majority,” in *Proceedings of the 19th Annual ACM Symposium on Theory of Computing, 1987, New York, New York, USA, 1987*, pp. 218–229.
  - [4] P. Paillier, “Public-key cryptosystems based on composite degree residuosity classes,” in *Advances in Cryptology - EUROCRYPT '99, International Conference on the Theory and Application of Cryptographic Techniques, Prague, Czech Republic, May 2-6, 1999, Proceeding*, 1999, pp. 223–238. [Online]. Available: [https://doi.org/10.1007/3-540-48910-X\\_16](https://doi.org/10.1007/3-540-48910-X_16)
  - [5] I. Damgård, M. Geisler, and M. Krøigaard, “Efficient and secure comparison for on-line auctions,” in *Australasian Conference on Information Security and Privacy*. Springer, 2007, pp. 416–430.
  - [6] C. Gentry, “Fully homomorphic encryption using ideal lattices,” in *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009, 2009*, pp. 169–178. [Online]. Available: <http://doi.acm.org/10.1145/1536414.1536440>
  - [7] J. Fan and F. Vercauteren, “Somewhat practical fully homomorphic encryption.” *IACR Cryptology ePrint Archive*, vol. 2012, p. 144, 2012.
  - [8] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, “(Leveled) fully homomorphic encryption without bootstrapping,” in *Proc. of ITCS*. ACM, 2012, pp. 309–325.
  - [9] C. Gentry, A. Sahai, and B. Waters, “Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based,” in *Advances in Cryptology—CRYPTO 2013*. Springer, 2013, pp. 75–92.
  - [10] P. Indyk and D. Woodruff, “Polylogarithmic private approximations and efficient matching,” in *Theory of Cryptography Conference*. Springer, 2006, pp. 245–264.
  - [11] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan, “Gazelle: A low latency framework for secure neural network inference,” in *27th USENIX Security Symposium*. USENIX Association, 2018.
  - [12] M. S. Riazi, C. Weinert, O. Tkachenko, E. M. Songhori, T. Schneider, and F. Koushanfar, “Chameleon: A hybrid secure computation framework for machine learning applications,” in *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*. ACM, 2018, pp. 707–721.
  - [13] A. Andoni, P. Indyk, and I. Razenshteyn, “Approximate nearest neighbor search in high dimensions,” *arXiv preprint arXiv:1806.09823*, 2018.
  - [14] J. Wang, H. T. Shen, J. Song, and J. Ji, “Hashing for similarity search: A survey,” *arXiv preprint arXiv:1408.2927*, 2014.
  - [15] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *arXiv preprint arXiv:1702.08734*, 2017.
  - [16] Y. A. Malkov and D. A. Yashunin, “Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs,” *IEEE transactions on pattern analysis and machine intelligence*, 2018.
  - [17] H. Jegou, M. Douze, and C. Schmid, “Product quantization for nearest neighbor search,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 1, pp. 117–128, 2011.
  - [18] A. Andoni, P. Indyk, T. Laarhoven, I. Razenshteyn, and L. Schmidt, “Practical and optimal lsh for angular distance,” in *Advances in Neural Information Processing Systems*, 2015, pp. 1225–1233.
  - [19] W. Microsoft Research, Redmond, “Simple Encrypted Arithmetic Library,” <http://sealcrypto.org>, 10 2018, SEAL 3.0.
  - [20] J. Doerner and A. Shelat, “Scaling ORAM for secure computation,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, 2017, pp. 523–535.
  - [21] A. Canteaut, S. Carpov, C. Fontaine, T. Lepoint, M. Naya-Plasencia, P. Paillier, and R. Sirdey, “Stream ciphers: A practical solution for efficient homomorphic-ciphertext compression,” in *Fast Software Encryption - 23rd International Conference, FSE 2016, Bochum, Germany, March 20-23, 2016, Revised Selected Papers*, 2016, pp. 313–333.
  - [22] D. G. Lowe *et al.*, “Object recognition from local scale-invariant features.” in *iccv*, vol. 99, no. 2, 1999, pp. 1150–1157.
  - [23] A. Babenko and V. Lempitsky, “Efficient indexing of billion-scale datasets of deep descriptors,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2055–2063.
  - [24] Z. Erkin, M. Franz, J. Guajardo, S. Katzenbeisser, I. Lagendijk, and T. Toft, “Privacy-preserving face recognition,” in *International Symposium on Privacy Enhancing Technologies Symposium*. Springer, 2009, pp. 235–253.
  - [25] A.-R. Sadeghi, T. Schneider, and I. Wehrenberg, “Efficient privacy-preserving face recognition,” in *International Conference on Information Security and Cryptology*. Springer, 2009, pp. 229–244.
  - [26] D. Evans, Y. Huang, J. Katz, and L. Malka, “Efficient privacy-preserving biometric identification,” in *Proceedings of the 17th conference Network and Distributed System Security Symposium, NDSS*, 2011.
  - [27] M. Barni, T. Bianchi, D. Catalano, M. Di Raimondo, R. Donida Labati, P. Failla, D. Fiore, R. Lazzeretti, V. Piuri, F. Scotti *et al.*, “Privacy-preserving fingerprint authentication,” in *Proceedings of the 12th ACM workshop on Multimedia and security*. ACM, 2010, pp. 231–240.
  - [28] D. Demmler, T. Schneider, and M. Zohner, “Aby-a framework for efficient mixed-protocol secure two-party computation.” in *NDSS*, 2015.
  - [29] P. Paillier, “Public-key cryptosystems based on composite degree residuosity classes,” in *International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 1999, pp. 223–238.
  - [30] H. Shaul, D. Feldman, and D. Rus, “Scalable secure computation of statistical functions with applications to  $k$ -nearest neighbors,” *arXiv preprint arXiv:1801.07301*, 2018.
  - [31] G. Asharov, S. Halevi, Y. Lindell, and T. Rabin, “Privacy-preserving search of similar patients in genomic data,” *Proceedings on Privacy Enhancing Technologies*, vol. 2018, no. 4, pp. 104–124, 2018.
  - [32] M. Aumüller, E. Bernhardsson, and A. Faithfull, “Ann-benchmarks: A benchmarking tool for approximate nearest neighbor algorithms,” in *International Conference on Similarity Search and Applications*. Springer, 2017, pp. 34–49.
  - [33] E. M. Songhori, S. U. Hussain, A.-R. Sadeghi, and F. Koushanfar, “Compacting privacy-preserving  $k$ -nearest neighbor search using logic synthesis,” in *Design Automation Conference (DAC), 2015 52nd ACM/EDAC/IEEE*. IEEE, 2015, pp. 1–6.
  - [34] M. S. Riazi, M. Javaheripi, S. U. Hussain, and F. Koushanfar, “MPCircuits: Optimized circuit generation for secure multi-party computation,” in *IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2019.
  - [35] D. Beaver, S. Micali, and P. Rogaway, “The round complexity of secure protocols,” in *STOC*, vol. 90, 1990, pp. 503–513.
  - [36] M. S. Riazi, B. Chen, A. Shrivastava, D. Wallach, and F. Koushanfar, “Sub-linear privacy-preserving near-neighbor search with untrusted server on large-scale datasets,” *arXiv preprint arXiv:1612.01835*, 2016.
  - [37] O. Goldreich and R. Ostrovsky, “Software protection and simulation on oblivious rams,” *Journal of the ACM (JACM)*, vol. 43, no. 3, pp. 431–473, 1996.
  - [38] D. Boneh, D. Mazieres, and R. A. Popa, “Remote oblivious storage: Making oblivious RAM practical,” 2011.
  - [39] K.-M. Chung, Z. Liu, and R. Pass, “Statistically-secure ORAM with  $\tilde{O}(\log^2 n)$  overhead,” in *International Conference on the Theory and Application of Cryptology and Information Security*. Springer, 2014, pp. 62–81.
  - [40] I. Damgård, S. Meldgaard, and J. B. Nielsen, “Perfectly secure oblivious RAM without random oracles,” in *Theory of Cryptography Conference*. Springer, 2011, pp. 144–163.
  - [41] O. Goldreich, “Towards a theory of software protection and simulation by oblivious RAMs,” in *Proceedings of the nineteenth annual ACM symposium on Theory of computing*. ACM, 1987, pp. 182–194.
  - [42] X. S. Wang, Y. Huang, T. H. Chan, A. Shelat, and E. Shi, “SCORAM: oblivious ram for secure computation,” in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2014, pp. 191–202.
  - [43] X. Wang, H. Chan, and E. Shi, “Circuit ORAM: On tightness of the goldreich-ostrovsky lower bound,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2015, pp. 850–861.
  - [44] S. Zahur, X. Wang, M. Raykova, A. Gascón, J. Doerner, D. Evans, and J. Katz, “Revisiting square-root ORAM: efficient random access in multi-party computation,” in *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2016, pp. 218–234.

[45] J. Doerner and A. Shelat, “Scaling oram for secure computation,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2017, pp. 523–535.

[46] V. Kolesnikov and T. Schneider, “Improved garbled circuit: Free XOR gates and applications,” in *Automata, Languages and Programming, 35th International Colloquium, ICALP 2008, Reykjavik, Iceland, July 7-11, 2008, Proceedings, Part II - Track B: Logic, Semantics, and Theory of Programming & Track C: Security and Cryptography Foundations*, 2008, pp. 486–498.

[47] S. Zahur, M. Rosulek, and D. Evans, “Two halves make a whole - reducing data transfer in garbled circuits using half gates,” in *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, April 26-30, 2015, Proceedings, Part II*, pp. 220–250.

[48] M. Bellare, V. T. Hoang, S. Keelveedhi, and P. Rogaway, “Efficient garbling from a fixed-key blockcipher,” in *2013 IEEE Symposium on Security and Privacy*. IEEE, 2013, pp. 478–492.

[49] S. Lloyd, “Least squares quantization in pcm,” *IEEE transactions on information theory*, vol. 28, no. 2, pp. 129–137, 1982.

[50] M. Blum, R. W. Floyd, V. R. Pratt, R. L. Rivest, and R. E. Tarjan, “Time bounds for selection,” *J. Comput. Syst. Sci.*, vol. 7, no. 4, pp. 448–461, 1973.

[51] M. Ajtai, J. Komlós, and E. Szemerédi, “An  $O(n \log n)$  sorting network,” in *Proceedings of the fifteenth annual ACM symposium on Theory of computing*. ACM, 1983, pp. 1–9.

[52] K. E. Batcher, “Sorting networks and their applications,” in *Proceedings of the April 30–May 2, 1968, spring joint computer conference*. ACM, 1968, pp. 307–314.

[53] P. Diaconis and D. Freedman, “Finite exchangeable sequences,” *The Annals of Probability*, pp. 745–764, 1980.

[54] M. Burkhart and X. Dimitropoulos, “Fast privacy-preserving top-k queries using secret sharing,” in *2010 Proceedings of 19th International Conference on Computer Communications and Networks*. IEEE, 2010, pp. 1–7.

[55] E. M. Kornaropoulos, C. Papamanthou, and R. Tamassia, “Data recovery on encrypted databases with k-nearest neighbor query leakage,” in *Data Recovery on Encrypted Databases with k-Nearest Neighbor Query Leakage*. IEEE, p. 0.

[56] Z. Brakerski, “Fully homomorphic encryption without modulus switching from classical gapsvp,” in *Advances in Cryptology - CRYPTO 2012 - 32nd Annual Cryptology Conference, Santa Barbara, CA, USA, August 19-23, 2012. Proceedings*, 2012, pp. 868–886.

[57] N. Gilboa and Y. Ishai, “Distributed point functions and their applications,” in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2014, pp. 640–658.

[58] E. Boyle, N. Gilboa, and Y. Ishai, “Function secret sharing,” in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2015, pp. 337–367.

[59] —, “Function secret sharing: Improvements and extensions,” in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016, pp. 1292–1303.

[60] D. Beaver, “Precomputing oblivious transfer,” in *Annual International Cryptology Conference*. Springer, 1995, pp. 97–109.

[61] V. Kolesnikov and T. Schneider, “Improved garbled circuit: Free xor gates and applications,” in *International Colloquium on Automata, Languages, and Programming*. Springer, 2008, pp. 486–498.

[62] S. Zahur, M. Rosulek, and D. Evans, “Two halves make a whole,” in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2015, pp. 220–250.

[63] P. Rindal, “libOTe: an efficient, portable, and easy to use Oblivious Transfer Library,” <https://github.com/osu-crypto/libOTe>.

[64] M. Aumüller, E. Bernhardtsson, and A. Faithfull, “Ann-benchmarks: A benchmarking tool for approximate nearest neighbor algorithms,” *Information Systems*, 2019.

[65] M. R. Albrecht, R. Player, and S. Scott, “On the concrete hardness of learning with errors,” *Journal of Mathematical Cryptology*, vol. 9, no. 3, pp. 169–203, 2015.

[66] J. Boyar and R. Peralta, “A small depth-16 circuit for the AES s-box,” in *Information Security and Privacy Research - 27th IFIP TC 11 Information Security and Privacy Conference, SEC 2012, Heraklion, Crete, Greece, June 4-6, 2012. Proceedings*, 2012, pp. 287–298.

[67] J. Doerner and A. Shelat, “Floram: The floram oblivious ram implementation for secure computation,” <https://gitlab.com/neucrypt/floram>.

[68] J. Doerner, “The absentminded crypto kit,” <https://bitbucket.org/jackdoerner/absentminded-crypto-kit>.

[69] D. J. Bernstein, “The salsa20 family of stream ciphers,” in *New Stream Cipher Designs - The eSTREAM Finalists*, 2008, pp. 84–97.

[70] —, “The chacha family of stream ciphers,” <https://cr.yp.to/chacha.html>.

[71] M. R. Albrecht, C. Rechberger, T. Schneider, T. Tiessen, and M. Zohner, “Ciphers for MPC and FHE,” in *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, April 26-30, 2015, Proceedings, Part I*, 2015, pp. 430–454.

[72] C. D. Canière and B. Preneel, “Trivium,” in *New Stream Cipher Designs - The eSTREAM Finalists*, 2008, pp. 244–266.

## APPENDIX

### A. Cost model

Here we describe the cost model we use to tune the hyperparameters. We focus on the clustering-based algorithm, since tuning the linear scan can be seen as an easy special case (all points in the stash etc.). We heavily use the notation introduced in Section V-F.

There are three main steps of the algorithm:

- 1) Compute closest  $u^i$  from the overall  $k_c^i$  centers for the  $i$ -th group for  $1 \leq i \leq T$ ;
- 2) Retrieve  $u_{\text{all}} = \sum_{i=1}^T u^i$  clusters from ORAM (of  $m$  points each);
- 3) Compute  $k_{\text{nn}}$  closest points from the union of the stash (of  $s$  points) and  $m \cdot u_{\text{all}}$  retrieved points.

1) *Number of AND gates:* As a proxy for the total cost, we use the number of AND gates in the circuits. For ORAM retrieval, we *do not* count AND gates necessary for the functional secret sharing.

- 1) Cost of computing closest centers of clusters:

$$\sum_{i=1}^T \left( k_c^i \cdot b_d + (k_c^i + u^i \cdot l^i) \cdot (2 \cdot (b_d - r_c) + b_{\text{cid}}) \right)$$

- 2) Cost of ORAM retrieval (modulo FSS):

$$u_{\text{all}} \cdot \left( 6 \cdot (1152 + m \cdot (d \cdot b_c + b_d + b_{\text{pid}})) + m \cdot (d+1) \cdot b_d \right)$$

- 3) Cost of computing closest points (the final answer):

$$(s + m \cdot u_{\text{all}}) \cdot b_d + \left( s + m \cdot u_{\text{all}} + k_{\text{nn}} \cdot (l_s + m \cdot u_{\text{all}}) \right) \times (2 \cdot (b_d - r_p) + b_{\text{pid}})$$

Our cost is defined as the sum of the three above expressions. Next, for completeness we list the formulae for the numbers of inputs and outputs for the server and client for all of the three parts. These quantities affect the communication, but we do not include them in the cost we optimize since they affect the computation time less than the number of AND gates.

- 1) Closest centers

- Server’s inputs:

$$\sum_{i=1}^T (k_c^i \cdot (b_d + b_{\text{cid}}) + u^i \cdot b_{\text{cid}})$$

- Client’s inputs:

$$\sum_{i=1}^T k_c^i \cdot (b_d + b_{cid})$$

- Client’s outputs:

$$\sum_{i=1}^T u^i \cdot b_{cid}$$

## 2) ORAM retrieval

- Server’s inputs:

$$u_{all} \cdot \left( b_{cid} + 128 + m \cdot (d \cdot b_c + b_d + b_{pid} + (d + 1) \cdot b_d + b_{pid}) \right)$$

- Client’s inputs:

$$u_{all} \cdot (b_{cid} + 128 + m \cdot (d \cdot b_c + b_d + b_{pid}))$$

- Client’s outputs:

$$u_{all} \cdot m \cdot ((d + 1) \cdot b_d + b_{pid})$$

## 3) Closest points

- Server’s inputs:

$$(s + u_{all} \cdot m) \cdot (b_d + b_{pid}) + k_{nn} \cdot b_{pid}$$

- Client’s inputs:

$$(s + u_{all} \cdot m) \cdot (b_d + b_{pid})$$

- Client’s outputs:

$$k_{nn} \cdot b_{pid}$$

## B. Hyperparameters for the evaluation

In Figure 5 and Figure 6, we summarize the parameters we use for both of our algorithms on each of the datasets. We find these parameters as approximate minimizers of our cost model from Section A.

## C. Stream Ciphers as PRF

In the original Floram construction [20], the PRF and the PRG used in the read-only process are chosen by the authors to be AES-128. Indeed, AES-128 is a block cipher that has been largely studied by the cryptographic community. The implementations of the scheme are highly optimized (less than 5000 non-free gates per block [66]) and its security is often used as a standard term of comparison. The implementation of Floram [67], [68] uses the optimized AES-128 and proposes two alternative symmetric encryption schemes: the streams Salsa20 [69] and its variant Chacha20 [70].

However, other symmetric ciphers can be used to obtain an efficient PRF/PRG. In particular, we looked for a PRF with low number of AND gates in order to decrease the communication between the parties when it is evaluated in GC (in the Free-XOR setting). Between the block ciphers, one of the most promising constructions is LowMC [71], which has a small number of AND gates per output bit. Between the stream

ciphers, instead, Trivium [72] and its variant Kreyvium [21] captured our attention. They are flexible in terms of input and output size, since there is no fixed block size to respect, and their evaluation is very efficient in terms of AND gates per output bit of stream.

Trivium belongs to the 2008 eSTREAM portfolio. It presents a simple construction, needing only 3 AND gates per bit of stream produced, plus  $3 \cdot 1152$  initialization AND gates executed once per stream. Trivium uses a secret key and an IV of size 80-bits each and achieves 80-bits of security. The scheme uses three registers, which are initialized with the key, the IV and some additional fixed bits. At each round, three temporary variables are computed by adding or multiplying some fixed elements in the register (9 XORs and 3 ANDs per round): at the end of each round, every register is rotated by 1 position, one element is discarded and one temporary value is appended. The first 1152 (this number is chosen for security reasons) rounds are the initialization rounds and they do not produce any stream. After the initialization, every round outputs one bit of stream, equal to the XOR of the three temporary values.

Kreyvium was presented in 2015 as a 128-bits secure variant of Trivium, as a solution particularly suited for homomorphic-ciphertext compression: the construction uses longer keys and IVs (128 bits each), 2 additional registers and a few additional XOR gates per round, but keeps the same amount of AND gates per bit of stream produced and for the initialization phase.

AES-128 needs about 5000 AND gates to produce 128 bits of stream, while Trivium and Kreyvium need  $3 \cdot (1152 + N)$  AND gates, where  $N$  is the size of the input/output of the stream. The difference is not impressive when the input blocks are of size 128, but the gap between the two ciphers increases when the size of inputs increases, since the stream cipher only needs 3 more AND gates per bit of input. For AES-128 the number of AND gates per bit remains constant (about 39 AND gates per output bit) while in Kreyvium it decreases to about 3 AND gates per bit of stream (see Table I).

The inputs we use in our construction have different sizes. For small datasets, every input is about 2.7 kB while for large datasets the inputs are about 5 or 6 kB. We compute 2 PRFs per input, so the actual number of AND gates in Table I should be doubled.

While our approach is more efficient in GC with respect to Floram, the plaintext evaluation of Kreyvium is slower than the (highly optimized) hardware implementation of AES. In order to mitigate this issue, we vertically batch 512 bits and we compute multiple streams in parallel (using AVX-512), so we are able to process several hundreds of Mega Bytes of information per second in single core.

## D. Optimal circuit for implicit top- $k$

Recall that our goal is, given  $n$  numbers each consisting of  $b$  bits, to find  $k$  smallest numbers in the following form: the output of a circuit is a binary vector with exactly  $k$  ones at



Parameter	Linear scan			Clustering		
	SIFT	Deep1B-1M	Deep1B-10M	SIFT	Deep1B-1M	Deep1B-10M
$l_s$	8334	8334	83	262	210	423
$b_c$	8	8	8	8	8	8
$r_p$	8	8	9	8	8	8

Fig. 5: (Near-)optimal hyperparameters that are used both by linear scan and the clustering-based algorithm.

Parameter	SIFT	Deep1B-1M	Deep1B-10M
$T$	4	5	6
$k_c^i$	50810 25603 9968 4227	44830 25867 11795 5607 2611	209727 107417 39132 14424 5796 2394
$m$	20	22	48
$u^i$	50 31 19 13	46 31 19 13 7	88 46 25 13 7 7
$s$	31412	25150	50649
$l^i$	458 270 178 84	458 270 178 84 84	924 458 178 93 84 84
$r_c$	5	5	5
$\alpha$	0.56	0.56	0.56

Fig. 6: (Near-)optimal hyperparameters that are specific to the clustering-based algorithm.

	128 bits	2.7 kB	6 kB
AES-128	5000 AND (39 AND/bit)	865000 AND (39.1 AND/bit)	1920000 AND (39.06 AND/bit)
Chacha20	20480 AND (160 AND/bit)	901120 AND (40.7 AND/bit)	1966080 AND (40 AND/bit)
Kreyvium	3840 AND (30 AND/bit)	69810 AND (3.15 AND/bit)	150912 AND (3.07 AND/bit)

TABLE I: Estimates on the number of AND gates for ciphers AES-128, Chacha20 and Kreyvium for different input sizes. The estimates for Chacha20 refer to a naive implementation of the scheme: we believe that the scheme would be more efficient in terms of non trivial gates in practice, but we did not found such optimal estimates in the literature. We do not report the number of AND gates for LowMC: they should be comparable to the estimates we have for Kreyvium for an optimal choice of the parameters.

the positions that correspond to the smallest elements. Such representation was used in [54] and [31].

Previously it was known how to achieve this in  $O(b^2n)$  gates as follows: we need to find a threshold  $y$  such that  $|\{i: x_i \leq y\}| = k$ , after that finding the result can be trivially done in  $O(bn)$  gates by comparing every number with  $y$ . We can find  $y$  using binary search, which takes  $b$  iterations, and for each iteration we compare every number with a current guess for  $y$ , which takes  $O(bn)$  gates, resulting in  $O(b^2n)$  gates overall.

Now we show how to improve this construction to the optimal  $O(bn)$  gates. Instead of running the full binary search for  $y$ , we will be computing it bit-by-bit starting from the most significant one. In order to do this, we maintain a binary  $a_i$  vector of “alive” elements of the list, initially  $a_i \equiv 1$ . To figure out  $i$ -th bit of  $y$ , we count how many alive elements of the

list have 0 as the  $i$ -th bit; let us denote this number by  $c_0$ . This counting can be done in  $O(n)$  gates using a binary tree of adders. Next we compare  $c_0$  with  $k$ : if  $k \leq c_0$ , then we zero out the entries of  $a$  for the elements of the list with the  $i$ -th bit being 1, otherwise, we zero out the entries with the  $i$ -th bit being 0 and subtract  $c_0$  from  $k$ . All of these operations can be implemented in  $O(n)$  gates, and there are  $b$  iterations in total. Overall, this circuit can be seen as a hybrid between radix sort and a randomized selection algorithm.