

Private Eyes: Zero-Leakage Iris Searchable Encryption

Julie Ha* Chloe Cachet† Luke Demarest‡ Sohaib Ahmad§ Benjamin Fuller¶

December 9, 2024

Abstract

This work introduces *Private Eyes*, the first zero-leakage biometric database. The only leakage of the system is unavoidable: 1) the log of the dataset size and 2) the fact that a query occurred. Private Eyes is built from oblivious symmetric searchable encryption. Approximate proximity queries are used: given a noisy reading of a biometric, the goal is to retrieve all stored records that are close enough according to a distance metric.

Private Eyes combines locality sensitive hashing or LSHs (Indyk and Motwani, STOC 1998) and oblivious maps which map keywords to values. One computes many LSHs of each record in the database and uses these hashes as keywords in the oblivious map with the matching biometric readings concatenated as the value. At search time with a noisy reading, one computes the LSHs and retrieves the disjunction of the resulting values from the map. The underlying oblivious map needs to answer disjunction queries efficiently.

We focus on the iris biometric which requires a large number of LSHs, approximately 1000. Boldyreva and Tang’s (PoPETS 2021) design yields a suitable map for a small number of LSHs (their application was in zero-leakage k -nearest-neighbor search).

Our solution is a zero-leakage disjunctive map designed for the setting when most clauses do not match any records. For the iris, on average at most 6% of LSHs match any stored value.

We evaluate using the ND-0405 dataset; this dataset has 356 irises suitable for testing. To scale our evaluation, we use a generative adversarial network to produce synthetic irises. Accurate statistics on sizes beyond available datasets is crucial to optimizing the cryptographic primitives. This tool may be of independent interest. For the largest tested parameters of a 5000 synthetic iris database, a search requires 18 rounds of communication and 25ms of parallel computation. Our scheme is implemented and open-sourced.

Searchable Encryption, Biometrics, Proximity Search

1 Introduction

Biometrics are collected into large databases for search [BBOH96, Dau14, Fou]. Learning stored biometric values enables attackers to break authentication and privacy for a user’s lifetime [GRGB⁺12, MCYJ18, AF20, VS11, HWKL18, SDDN19]. To reduce this risk, this article develops new searchable encryption techniques for biometric databases [SWP00, CGKO11]. See previous reviews of searchable encryption [BHJP14, FVY⁺17, KKM⁺22].

Consider a data owner outsourcing a database \mathcal{DB} to an honest but curious server that may learn information called *leakage*. Prior work exploits such leakage to reveal sensitive information about the database or queries [IKK12, CGPR15, KKNO16, WLD⁺17, GSB⁺17, GLMP18, KPT19, MT19, KE19, KPT20a, FMC⁺20, FP22, GPP23, KKM⁺22]. Since biometrics cannot be replaced or revoked, we focus on a **zero-leakage system**. **A zero-leakage system reveals only unavoidable information: 1) that a query occurs and 2) $|\mathcal{DB}|$ (which we pad to a power of 2)**. Our focus is on an efficient combination of zero-leakage primitives for this task. Components can be replaced with higher-leakage alternatives for efficiency gains.

Most biometric search systems phrase proximity queries as a large disjunction. For actual biometrics, 1) this disjunction has hundreds or thousands of terms and 2) most clauses will not match any stored record. Our parameter analysis shows that these specificities of biometrics result in critical inefficiencies in previous designs. We propose a new design that avoids these pitfalls.

*Boston University. Email hajulie@bu.edu

†University of Connecticut. Email chloe.cachet@uconn.edu

‡University of Connecticut. Email luke.h.demarest@gmail.com

§University of Connecticut. Email sohaib.ahmad@uconn.edu

¶University of Connecticut. Email benjamin.fuller@uconn.edu

Our Contributions We present `PrivateEyes`, the first zero-leakage proximity search system for the iris. It is enabled by the following contributions:

1. A two-stage design that 1) finds the non-null clauses (i.e. clauses that match something in the database), 2) obviously searches these non-null clauses *only*. The first stage can be built with private set intersection (PSI) [FNP04]. The second stage is built using oblivious tree traversal [WNL⁺14].
2. A prototype implementation with evaluation on random data up to 25000 records and synthetic iris data up to 5000 records [pri24]. On synthetic data, search takes at most 25ms of parallel computation and 18 rounds of communication. 92% of the computation and 14 rounds are in the second stage, the focus of this work.

To scale beyond the sizes of available iris datasets we also introduce a synthetic iris generation tool that may be of independent interest.

Organization The rest of this paper is organized as follows: Section 2 reviews our design including relevant prior work, Section 3 introduces preliminaries, Section 4 presents the details of our system, Section 5 presents the datasets, Section 6 describes our implementation, Section 7 our evaluation, and Section 8 concludes. Appendix A describes the architecture that generates synthetic irises.

2 Design Overview

A database is a list of biometrics $\mathcal{DB} = w_1, \dots, w_\ell$ where each $w_i \in \{0, 1\}^n$. The goal of a biometric database is, given some biometric query w^* , to find all values $w_i \in \mathcal{DB}$ that are close enough to w^* . For the Hamming metric \mathcal{D} and distance threshold t , that means to find all w_i such that $\mathcal{D}(w_i, w^*) \leq t$.¹

Like prior work [KIK12, FWG⁺16, WYLH14, LPW⁺20], we combine locality-sensitive hashes (LSHs) [IM98] with an encrypted map.²

An LSH maps *near* items to the same value more frequently than it maps *far* items to the same value. Let \mathcal{H} be a family of LSHs then

$$\Pr_{\text{LSH} \leftarrow \mathcal{H}} [\text{LSH}(w_i) = \text{LSH}(w^*) | w_i, w^* \text{ are near}] \geq 1 - p_1,$$

$$\Pr_{\text{LSH} \leftarrow \mathcal{H}} [\text{LSH}(w_j) = \text{LSH}(w^*) | w_j, w^* \text{ are far}] \leq 1 - p_2.$$

where $p_1 < p_2$.

Maps associate keywords to values and are used to build inverted indices. For a map M , we use the notation $M[\text{keyw}] = \text{value}$ to denote that `keyw` is associated with `value`.

For a database of size ℓ , parameter $\beta \in \mathbb{Z}^+$, LSH family \mathcal{H} , and maps M_1, \dots, M_β , one can achieve proximity search as follows:

1. Sample β LSHs, one per map: $\text{LSH}_1, \dots, \text{LSH}_\beta \leftarrow \mathcal{H}$.
2. In each map, **associate all records with the same LSH output to that LSH output**: for $j = 1, \dots, \beta$, set

$$M_j[v] = \{w_i | \text{LSH}_j(w_i) = v\}.$$
³

3. To search for value w^* , **retrieve all records with at least one LSH output in common with w^*** :

- (a) Compute $\text{LSH}_1(w^*), \dots, \text{LSH}_\beta(w^*)$.
- (b) Retrieve $\cup_{j=1}^\beta M_j[\text{LSH}_j(w^*)]$.

¹This functionality differs from k -nearest neighbors where the goal is to retrieve the k closest records [BT21]. There have been leakage abuse attacks against k -nearest neighbor systems that reveal access pattern [KPT19, KPT20b, LMWY20] and resulting systems [CCD⁺20]. These attacks do not apply to our leakage profile.

²We won't discuss works that use encrypted maps but require work proportional to the total number of close points, making them impractical for biometrics [LWW⁺10, WMT⁺13, BC14].

³If multiple records share the same LSH value our implementation concatenates the matching values. This allows us to handle a constant number of values associated to each keyword. This condition is satisfied for the accuracy regimes discussed in this work.

Queries are disjunctions. Boldyreva and Tang [BT21] constructed a zero-leakage encrypted map scheme called an *oblivious map with encryption* or OMapE. Each clause is submitted to the relevant OMapE, the results are concatenated as in Step 3b above.

Due to their large noise, **biometrics require sampling hundreds or thousands of LSHs** to achieve reasonable accuracy (see analysis in Section 3.1 and Section 7). At the same time, **few of these LSHs match anything** in the corresponding map. Constructions frequently perform heavy oblivious operations to hide the null value.

Our design separates the tasks of identifying which values are null and finding the associated values. That is, we first try to find a small number, called δ , of LSHs values that exist in some map, and then query exactly those δ maps in a way that hides which maps are being queried. For the above design to be successful, one needs to demonstrate:

1. **Obliviousness** One can hide the queried δ maps,
2. **Accuracy** High accuracy with $\delta < \beta$, and
3. **Speed** The approach is faster.

We now provide a more detailed description of the approach, shown visually in Figure 1.

Oblivious Membership Check An object to check which LSHs’ outputs, $\text{LSH}_j(w^*)$, have matches. For an encrypted stored set X the *oblivious membership check* or OMC takes in a set W and returns a set W' of size δ . If $|W \cap X| \leq \delta$ then W' contains $W \cap X$ and $\delta - |W \cap X|$ dummy values. If $|W \cap X| > \delta$ then W' contains δ randomly selected values from $W \cap X$. In our system, X is the set of all LSHs values $X = \{(j, \text{LSH}_j(w_i))\}_{i,j}$ and $W = \{(j, v) | \text{LSH}_j(w^*) = v\}_j$ for some query w^* .

We build OMC using private set intersection and pseudorandom permutations. We benchmark this design using VolePSI [RS21]; the resulting implementation takes *2ms* and 4 rounds of communication for our largest tested parameters.

For each parameter set, we manually find a size δ where if one only searches for δ items in the second stage, there is only a small degradation in accuracy.

Disjunctive Oblivious Map An object that searches for the disjunction of exactly δ items. These δ values form a set of LSHs that exist in the map⁴. Using an oblivious data structure with a constant number of queries yields a zero-leakage solution. This object has the same functionality as an oblivious map that takes multiple clauses but the fact that all clauses are presented together is crucial for security. We call this object a DOMapE for *disjunctive oblivious map with encryption*. Our focus is on designing a DOMapE.

We propose a DOMapE based on oblivious tree traversal, building on the design principles of Wang et al. [WNL⁺14]. In our approach, each map corresponds to a tree, so one always performs δ tree traversals. We use oblivious RAMs to store tree nodes. The oblivious RAMs are organized to minimize nodes that are stored together while ensuring that δ tree traversals result in no leakage.

To build the encrypted database (Enc stage in figure 1) the client computes the set of all LSH values using the OMC (with their index j). They also create the maps associating LSH values to the corresponding biometric values and build the corresponding DOMapE. To search for query w^* (Find stage in figure 1), the client computes the corresponding LSH values and uses OMC.Find to find a subset of those values present in the maps. DOMapE.Find is then used to retrieve the corresponding non-null values from the maps. The inputs to both OMC.Find and DOMapE.Find are of constant size, yielding a zero-leakage solution.

We implemented and analyzed this construction on datasets of up to 25000 records. Since available iris datasets are not this large, we created larger synthetic iris datasets up to 5000 records using generative adversarial networks or GANs [CWD⁺18, GPAM⁺20] and random data to 25K records. For all tested parameters we have $\delta < 75$ (in contrast to a number of LSHs $\beta > 600$), so a server could reasonably process the ORAM requests at each tree level in parallel. Thus, we report on both parallel and sequential time.

Comparison with prior work Cachet et al. [ACD⁺22] proposed two non-interactive iris proximity search schemes based on inner-product encryption. Both of their constructions have more leakage than our system. The first one leaks the distance between all returned points and the query. The other leaks whether returned records are the same

⁴Except for dummy values added when the total number of matches is less than δ .

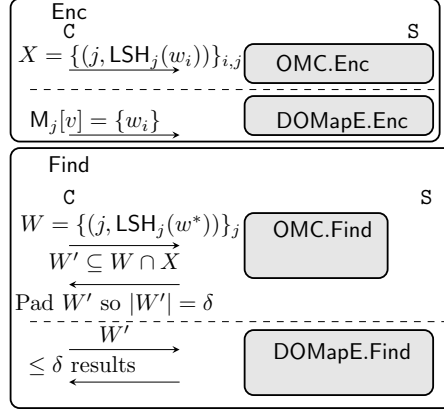


Figure 1: System Overview Composing OMC and DOMapE. In Find one checks which LSH values exist in the database using the OMC, finding at most δ candidates. If needed OMC pads this set of candidates with dummy values to reach $|W'| = \delta$. DOMapE is designed to have no leakage if the size of the disjunction is constant (δ). See Construction 2 for a formal description.

distance from the query. For the solution with more leakage, search took 4 minutes on a dataset of size 356. For the solution with less leakage, the search took 1 hour [CAD⁺23].

Barni et al. [BBC⁺10] and Blanton et al. [BG11] evaluated their fingerprints identification systems on small datasets of 320 readings. They respectively achieved search times of 16 and 0.45 seconds for the server. Their protocols use Garbled circuits and only a single round of communication. Blanton et al. [BG11] also proposed an iris identification system that compares two 2048 bits iris readings in 0.15s. For a database of size 356 (our smallest setting) this would amount to 53s. For face recognition, SciFi [OPJM10] online search runtime is linear in the size of the database: for 100 faces representations search takes 31s. Erkin et al. [EFG⁺09]’s system takes roughly 40s to search over a dataset of size 320. As a reminder on 5000 irises our search is 25ms with 18 rounds. For a network delay of 100ms, this corresponds to an overall time of $\approx 2s$.

Recent work [UCK⁺21, GRS22, CFR23] extends PSI to the setting where one considers items a match if their distance is small. A naive use of fuzzy PSI tells the client if the biometric exists in the database but not which item it matches. This issue can be solved using labeled PSI, [CHLR18] which associates a label with each set value $x \in X$. This is the approach used by Uzun et al. [UCK⁺21]. Uzun et al. [UCK⁺21] evaluate the face biometric. At a technical level, their fuzzy PSI is similar to our approach; they store LSH outputs in the set X and encrypt each item before sending it to the server. However, they perform a threshold version of labeled PSI, only returning a record if there are enough matches. One match is required in our approach. Their approach is bandwidth efficient but requires fully homomorphic encryption [Gen09] to perform the more complex matching. Uzun et al’s scheme does not allow the client to learn about non-matching records which is not a goal of our work. Their scheme uses $\beta = 64$ by averaging multiple biometric readings, such techniques can be applied in our setting as well. Averaging readings is not standard in the biometric literature so we do not do it.

More work is needed to scale zero-leakage biometric databases. Our cryptographic storage overhead is around 5000 for all tested parameters. The required number of LSHs grows with the database size. Scaling to 10^6 irises requires 93000 LSHs (using linear interpolation over the number of required LSHs for synthetic data with an $r^2 = .88$) for an overall storage of 93 billion LSH outputs (roughly 500TB for the cryptographic object). In our current testing, we load the whole cryptographic object in memory with each query, which is not possible for larger sizes.

3 Preliminaries

Let λ be the security parameter, we use $\text{poly}(\lambda)$ and $\text{negl}(\lambda)$ to denote unspecified functions that are polynomial and negligible in λ , respectively. All definitions are indexed by λ but this indexing is omitted for notational clarity. For some $n \in \mathbb{N}$, $[n]$ denotes the set $\{1, \dots, n\}$. Let $x \leftarrow S$ denote sampling x uniformly at random from the finite set S . Let \perp_1, \perp_2, \dots be a sequence of distinguished unique symbols. These symbols are allowed inputs to algorithms. For a map M , let $M.\text{Keywords}$ output the set of all stored keywords.

For interactive protocol Prot between a client \mathbf{C} and a server \mathbf{S} we use notation $\begin{pmatrix} o_{\mathbf{C}} \\ o_{\mathbf{S}} \end{pmatrix} \leftarrow \text{Prot} \begin{pmatrix} i_{\mathbf{C}} \\ i_{\mathbf{S}} \end{pmatrix}$ with $i_{\mathbf{C}}, o_{\mathbf{C}}, i_{\mathbf{S}}, o_{\mathbf{S}}$ denoting the client’s and the server’s inputs and outputs respectively. Protocols are written from the perspective of the client with underlying interactive algorithms indicating the server’s role.

Hamming distance is defined as the distance between the bit vectors x and y of length n , $\mathcal{D}(x, y) = |\{i \mid x_i \neq y_i\}|$. The fractional Hamming distance is $\mathcal{D}(x, y)/n$.

Definition 1 (Locality-sensitive Hashing). *Let $t \in \mathbb{N}$, $c > 1$ and $p_1, p_2 \in [0, 1]$. \mathcal{H} defines a (t, ct, p_1, p_2) -sensitive hash family if for any $x, y \in \{0, 1\}^n$, we have:*

1. If $\mathcal{D}(x, y) \leq t$, $\Pr_{\text{LSH} \leftarrow \mathcal{H}} [\text{LSH}(x) = \text{LSH}(y)] \geq 1 - p_1$ and
2. If $\mathcal{D}(x, y) \geq ct$, $\Pr_{\text{LSH} \leftarrow \mathcal{H}} [\text{LSH}(x) = \text{LSH}(y)] \leq 1 - p_2$.

For x, y if $\mathcal{D}(x, y) \leq t$ they are called *near*, if $\mathcal{D}(x, y) \geq ct$ they are called *far*.

We use selection of a single random bit as our LSH. For two values x, y the probability this bit will be the same is $1 - \mathcal{D}(x, y)/n$. That is $p_1 \geq t/n$ while $p_2 \leq ct/n$. The error rates p_1, p_2 of an LSH can be increased by randomly sampling several LSHs and checking that they all match, an α -AND. For α -AND composition, a (t, ct, p_1, p_2) -LSH yields a $(t, ct, p_1^\alpha, p_2^\alpha)$ -LSH. For our LSH, this corresponds to randomly selecting α (with replacement) bits of the input. Similarly, p_1, p_2 can be decreased by randomly sampling several LSHs and checking that at least one of them matches, a β -OR. For β -OR composition, a (t, ct, p_1, p_2) -LSH yields a $(t, ct, 1 - (1 - p_1)^\beta, 1 - (1 - p_2)^\beta)$ -LSH.

3.1 The need for many LSHs in biometric proximity search

We focus on the iris biometric and use the ND-0405 dataset [PSO⁺09, BF16]. The average distance between same irises is $t/n \approx .21$, using a state-of-the-art feature extractor [AF19]. Section 5 presents the datasets used in this work. Our discussion applies to other biometrics with substantive noise such as the face⁵.

If one only uses the AND of α LSHs and considers it a match when at least one LSH out of β corresponds, this LSH can be seen as the β -OR of the α -AND of LSHs where

$$p'_1 = 1 - (1 - p_1^\alpha)^\beta \quad \text{and} \quad p'_2 = 1 - (1 - p_2^\alpha)^\beta.$$

Let w'_i be a noisy reading of w_i . When using w'_i as a search query, a true accept is when w_i is returned and the true accept rate (TAR) is the fraction of queries where this happens. The fraction of false accepts (FFA) is the fraction of $\mathcal{DB} \setminus \{w_i\}$ that is returned on average.

If one assumes that $p_1 = t/n = .21$ and $p_2 = .5$ and all items have these average distances, then TAR is $1 - p'_1$ and FFA is $1 - p'_2$. Under this assumption, achieving a TAR of .95 ($p'_1 = .05$) and an FFA of .01 ($p'_2 = .99$) requires a number of LSHs $65 \leq \beta \leq 80$ for the minimum $\alpha = 13$. For a dataset of size 10^6 if one seeks at most 100 false accepts (that is, FFA of 10^{-4}) this requires $680 \leq \beta \leq 835$ at the minimum $\alpha = 23$.

However, even though the mean distance between readings of the same biometric is $t/n = .21$ there is substantial variance in this distance (see Figure 11(a)), requiring β to be larger as we show in Table 2.

We call an LSH match *good* if it ensures the query results in a true accept and *bad* otherwise. For the ND-0405 dataset,⁶ a histogram of the number of good and bad LSH matches is in Figure 2. The average number of total LSH matches is 23.4.

3.2 Cryptographic Definitions

This work relies on *oblivious RAM* (ORAM) [GO96, Gol87] to achieve zero-leakage. In our constructions, we consider static datasets and ignore write queries. Our ORAM definition reflects this choice. As we discuss in Section 6, this definition is also satisfied by private information retrieval schemes (with appropriate encryption). We discuss other considerations for dynamic data at the end of Section 6.

⁵Deng et al. [DGXZ19, Figure 6] show analogous statistics for the face.

⁶This uses the following experiment:

1. Storage of a single feature extracted reading for the eye for each of the 356 persons. Sample $\beta = 225$ LSHs of size $\alpha = 15$.
2. Let w'_1, \dots, w'_{356} be the second stored template in the ND-0405 dataset.
3. Search for each record w'_i . Record the number of good and bad LSH matches.

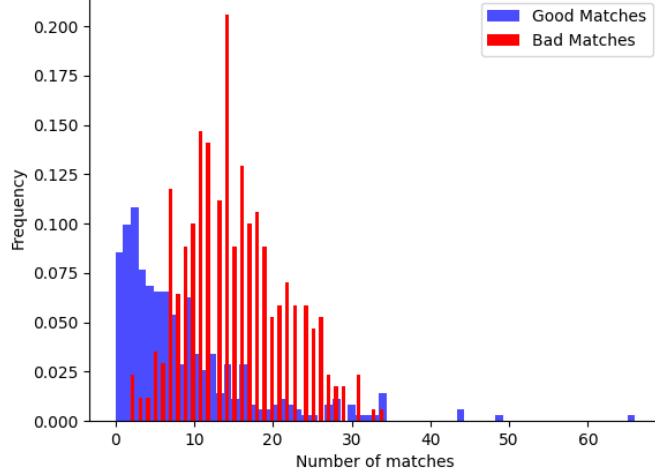


Figure 2: The number of maps returning a value when searching similar values with $\alpha = 15$, $\beta = 225$ trees and 356 records from the ND-0405 dataset.

Definition 2 (Oblivious RAM). *An Oblivious RAM (ORAM) scheme is two protocols, Setup and Access:*

$$\begin{pmatrix} \sigma \\ \text{EM} \end{pmatrix} \leftarrow \text{Setup} \left(1^\lambda, \text{Mem} \right), \begin{pmatrix} v, \sigma' \\ \text{EM}' \end{pmatrix} \leftarrow \text{Access} \left(\begin{pmatrix} \sigma \\ \text{EM} \end{pmatrix}, i \right).$$

Correctness *Consider the following correctness experiment:*

1. *An adversary \mathcal{A} chooses memory Mem.*

2. *Consider $\begin{pmatrix} \sigma_0 \\ \text{EM}_0 \end{pmatrix} \leftarrow \text{Setup} \left(1^\lambda, \text{Mem} \right)$.*

3. *For $1 \leq i \leq q$:*

(a) *Run $y_i \leftarrow \mathcal{A}(\text{ts}_{i-1})$.*

(b) *Run $\begin{pmatrix} v_i, \sigma_i \\ \text{EM}_i \end{pmatrix} \leftarrow \text{Access} \left(\begin{pmatrix} \sigma_{i-1} \\ \text{EM}_{i-1} \end{pmatrix}, y_i \right)$.*

The adversary wins if for some $i, v_i \neq \text{Mem}[y_i]$. The ORAM scheme is correct if the probability of \mathcal{A} winning the game is $\text{negl}(\lambda)$.

Security *An ORAM scheme is secure in the semi-honest model if for any PPT adversary \mathcal{A} , there exists a PPT simulator Sim such that for any PPT distinguisher D we have*

$$\left| \Pr[D(\text{Real}_{\mathcal{A},q}(1^\lambda)) = 1] - \Pr[D(\text{Ideal}_{\mathcal{A},\text{Sim},q}(1^\lambda)) = 1] \right| \leq \text{negl}(\lambda)$$

with $\text{Real}_{\mathcal{A},q}$ and $\text{Ideal}_{\mathcal{A},\text{Sim},q}$ as described in Figure 3.

The above is an adaptive simulation definition of ORAM [GMP16], all of our proofs work naturally for the standard non-adaptive definition. We define generic oblivious searchable encryption (OSE) and in the rest of the paper, will use specific variants of it.

Definition 3 (Oblivious searchable encryption). *Let \mathcal{M} denote the records space, \mathcal{Q} denote the query space and \mathcal{R} denote the result space. Let $\text{DB} \subseteq \mathcal{M}$ be a database and $y \in \mathcal{Q}$ be a query. For string param, let the triple of protocols*

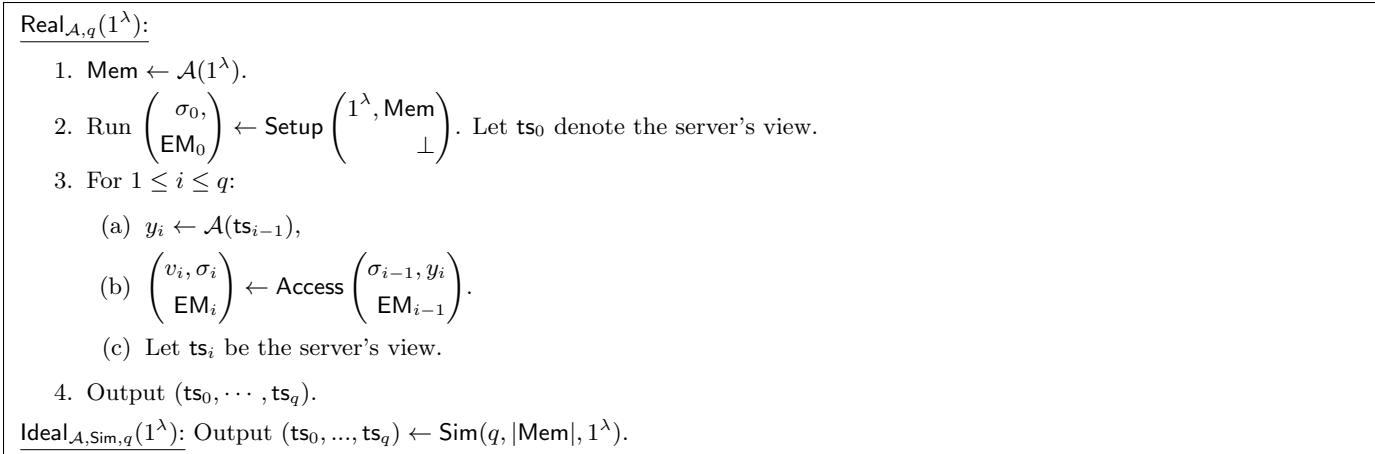


Figure 3: Definition of Ideal and Real for ORAM security.

OSE = (Setup, Enc, Find) have the following format:

$$\begin{pmatrix} \text{sk} \\ \text{pp} \end{pmatrix} \leftarrow \text{Setup} \begin{pmatrix} 1^\lambda, \text{param} \\ 1^\lambda, \text{param} \end{pmatrix},$$

$$\begin{pmatrix} I_c, \\ I_s \end{pmatrix} \leftarrow \text{Enc} \begin{pmatrix} \text{sk}, \mathcal{DB} \\ \text{pp} \end{pmatrix}, \begin{pmatrix} J, I'_c \\ I'_s \end{pmatrix} \leftarrow \text{Find} \begin{pmatrix} \text{sk}, y, I_c, \\ \text{pp}, I_s \end{pmatrix}.$$

OSE is an oblivious searchable encryption if the following hold:

Correctness: The set J is the “same” as the result of the query. The formal definition varies per OSE variant we consider and is defined later.

Security: Let $q = \text{poly}(\lambda)$, $\mathcal{L}_{\text{OSE}} = \{\mathcal{L}_{\text{Enc}}, \mathcal{L}_{\text{Find}} = \perp\}$ be the leakage profile of OSE's algorithms. For any PPT adversary \mathcal{A} , there exists a simulator Sim such that for any PPT distinguisher D we have

$$|\Pr[D(\text{Real}_{\mathcal{A},q}(1^\lambda)) = 1] - \Pr[D(\text{Ideal}_{\mathcal{A},\text{Sim},q}(1^\lambda)) = 1]| \leq \text{negl}(\lambda)$$

with $\text{Real}_{\mathcal{A},q}$ and $\text{Ideal}_{\mathcal{A},\text{Sim},q}$ as described in Figure 4.

Our goal is to build an OSE scheme for proximity queries, we define this particular variant of OSE as follows:

Definition 4 (Oblivious Proximity Search). Consider Definition 3 with the following specificities. Let $\mathcal{M} = \mathcal{Q} = \mathcal{R} = \{0, 1\}^n$ and $\text{param} = t$. Consider $\mathcal{DB} = w_1, \dots, w_\ell$ where each $w_i \in \mathcal{M}$.

(ϵ, t)-Approximate Correctness: For all $\mathcal{DB}, y \in \mathcal{Q}$ define $J_{\mathcal{DB}, \text{near}, y} := \{w_i | \mathcal{D}(w_i, y) \leq t\}$. Let $q = \text{poly}(\lambda)$ and $\epsilon > 0$. For all \mathcal{DB} and all y_1, \dots, y_q define:

$$\begin{pmatrix} \text{sk} \\ \text{pp} \end{pmatrix} \leftarrow \text{Setup} \begin{pmatrix} 1^\lambda \\ 1^\lambda \end{pmatrix},$$

$$\begin{pmatrix} I_c^1, \\ I_s^1 \end{pmatrix} \leftarrow \text{Enc} \begin{pmatrix} \text{sk}, \mathcal{DB} \\ \text{pp} \end{pmatrix}, \begin{pmatrix} J_j, I_{c,j+1} \\ I_{s,j+1} \end{pmatrix} \leftarrow \text{Find} \begin{pmatrix} \text{sk}, y_j, I_{c,j} \\ \text{pp}, I_{s,j} \end{pmatrix}$$

OSE is ϵ -approximately correct if $\forall 1 \leq j \leq q$ for all \mathcal{DB}

$$\Pr [J^j \supseteq J_{\mathcal{DB}, \text{near}, y}] \geq 1 - \epsilon.$$

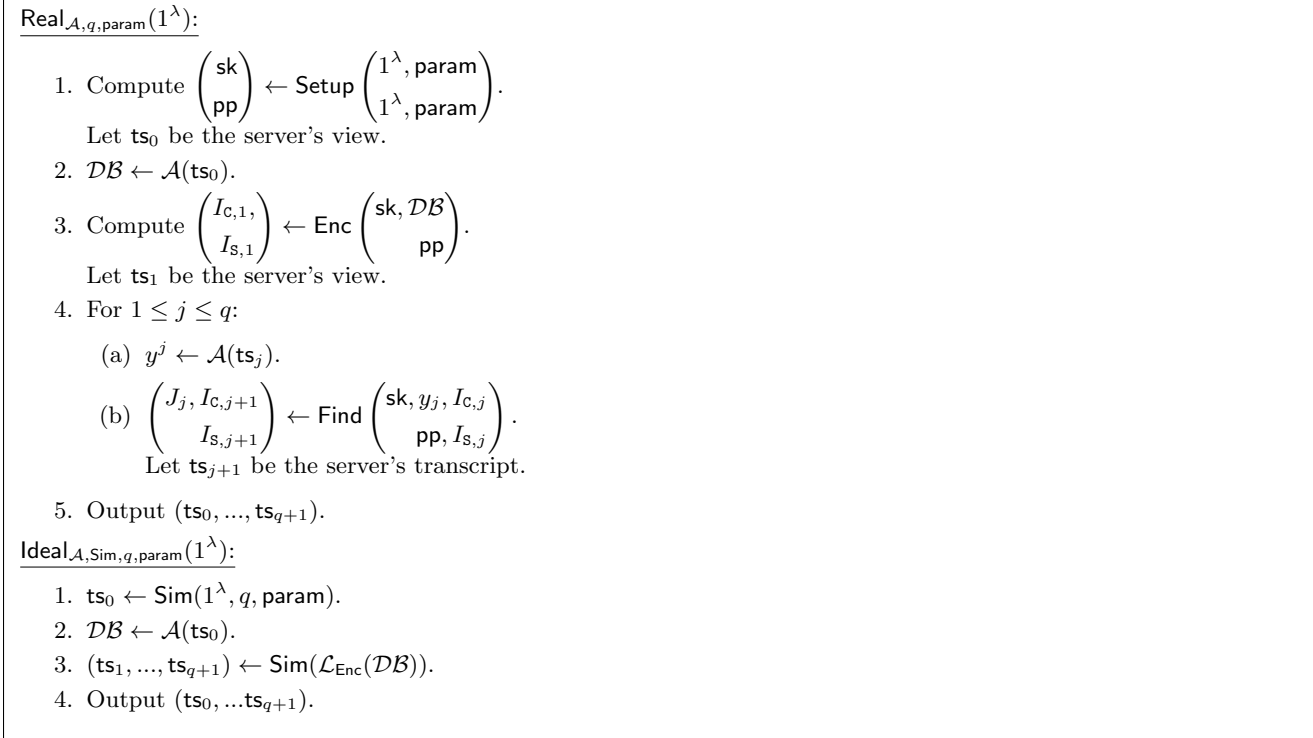


Figure 4: Definition of Ideal and Real for OSE security.

Definition 4 doesn't limit the number of false matches. Furthermore, in Section 4, we never show that our construction satisfies approximate correctness. Instead, we evaluate approximate correctness using data in Section 7.

Our scheme first finds a list of candidate LSH matches, and then uses an appropriate oblivious map to find the relevant records using the candidate LSH matches. The first stage is called *oblivious membership checking* or OMC. An OMC can be built from private set intersection (PSI), client storage, and full set retrieval (see Section 4.4). We are **not** offering constructions of OMC as a technical contribution. We benchmark separately using PSI, see discussion in Section 7. In our full implementation we use a local Bloom filter to simplify evaluation.

OMC only handles sets, that is, a collection of values *without repeats*. In our search system, these values are LSH outputs. It is possible for two distinct LSHs to have the same output. To avoid this, we prepend the LSH id to each LSH output value. For LSH j , the corresponding values to use would then be $\{j \parallel \text{LSH}_j(x)\}$. We define OMC as a variant of OSE:

Definition 5 (Oblivious Membership Check). *Let OMC = (Enc, Find) be a pair with stored set size ρ , query size γ , and result size δ , abbreviated ρ -ssize, γ -qsize, and δ -rsize. Consider Definition 3 with the following specificities:*

- Let $\mathcal{M} = \mathcal{Q} = \mathcal{R}$ and $\text{param} = (\rho, \gamma, \beta)$.
- Consider $X \subseteq \mathcal{M}$, such that $|X| = \gamma$, and $Y \subseteq \mathcal{Q}$, such that $|Y| = \rho$. Set $\mathcal{DB} = X$ and query $y = Y$.

Correctness: We use $\perp_1, \dots, \perp_\beta$ to denote a sequence of unique symbols that cannot appear in X or Y . For all $X, |X| = \gamma$ and $Y, |Y| = \beta$, let

$$\begin{pmatrix} EC \\ ES \end{pmatrix} \leftarrow \text{Enc} \begin{pmatrix} 1^\lambda, X \\ 1^\lambda \end{pmatrix}, \begin{pmatrix} I \\ \perp \end{pmatrix} \leftarrow \text{Find} \begin{pmatrix} EC, Y \\ ES \end{pmatrix}.$$

Then $|I| = \delta$ and for all $i \in I$ such that $\forall j, i \neq \perp_j$ it holds that $\Pr[i \in X \cap Y] \geq 1 - \text{negl}(\lambda)$.

Finally, we define the second stage of our system:

Parameter	Meaning
ℓ	Database size
β	nmaps , Total number of trees, maps, LSHs, query size to OMC
δ	qsize , Result size from OMC, Number of trees to traverse/result size in DOMapE
μ	msize , \log_2 maximum number of keywords in each map
γ	Server set size in OMC, in composed protocol $\gamma = 2^\mu \cdot \beta$

Table 1: Parameters and their usage across the different schemes.

Definition 6 (Disjunctive Oblivious Map with Encryption). *Let $\mathit{param} = (\beta, \delta)$, such that $\beta, \delta \in \mathbb{N}$ and $\delta \leq \beta$ and let $\mu \in \mathbb{N}$. Let $\text{DOMapE} = (\text{Setup}, \text{Enc}, \text{Find})$ be a triple with β **nmaps**, μ **map size**, and δ **query size**, abbreviated $\beta - \text{nmaps}$, $\mu - \text{msize}$, and $\delta - \text{qsize}$. Then DOMapE is a disjunctive oblivious map with encryption if it satisfies Definition 3 with the following correctness guarantee.*

Correctness: *Let $\mathcal{M} = \{M_i \mid M_i : \mathcal{Q} \leftarrow \mathcal{R}\}$, where M_i denotes a map such that for $1 \leq i \leq \beta$, $|M_i.\text{Keywords}| \leq 2^\mu$. Set $\mathcal{DB} = M_1, \dots, M_\beta$. Let $\epsilon > 0$, $q, \beta, \delta = \text{poly}(\lambda)$ and $\delta \leq \beta$. Let $\mathit{param} = (\beta, \mu, \delta)$. Fix some $(\{M_i\}_{i=1}^\beta, \{y^j \in (\mathcal{X} \times [1, \ell]^\delta)\}_{j=1}^q})$ and define for $1 \leq j \leq q$:*

$$\begin{aligned} \begin{pmatrix} \text{sk} \\ \text{pp} \end{pmatrix} &\leftarrow \text{Setup} \left(\begin{pmatrix} 1^\lambda, \beta, \delta \\ 1^\lambda \end{pmatrix} \right), \\ \begin{pmatrix} \sigma_1, \\ \text{EM}_1 \end{pmatrix} &\leftarrow \text{Enc} \left(\begin{pmatrix} \text{sk}, \mathcal{DB} \\ \text{pp} \end{pmatrix} \right), \begin{pmatrix} r_j, \sigma_{j+1}, \\ \text{EM}_{j+1} \end{pmatrix} &\leftarrow \text{Find} \left(\begin{pmatrix} \text{sk}, \sigma_j, y_j \\ \text{pp}, \text{EM}_j \end{pmatrix} \right). \end{aligned}$$

DOMapE is correct if there exists a set $\mathcal{I} \subseteq [2^\mu]$ where $|\mathcal{I}| \leq \delta$ such that :

$$\Pr \left[(\cup_i r_i^j) \setminus \emptyset \subseteq \cup_{i \in \mathcal{I}} M_{k_i} [x_i^j] \right] \geq 1 - \text{negl}(\lambda).$$

4 Oblivious Proximity Search for Biometrics

This section presents our technical solutions, focusing on the design of DOMapE . We describe possible constructions of OMC in Section 4.4. The most relevant related work is by Boldyreva and Tang [BT21], whose construction is for the approximate k -nearest neighbors search problem. While Boldyreva and Tang discuss two ways of implementing OMapE , one using a tree and the other using a skip list [Pug90], we present a tree based construction. In this work, we only consider static data. For static data, B-trees and skip lists are equivalent data structures [LN⁺96]. However, updates and the resulting performance differ.

4.1 Overview of DOMapE design

Recall the unprotected solution for proximity search from the Introduction:

1. Sample β LSHs, $\text{LSH}_1, \dots, \text{LSH}_\beta \leftarrow \mathcal{H}$.
2. For $j = 1, \dots, \beta$, set

$$M_j[\text{keyw}] = \{w_i \mid \text{LSH}_j(w_i) = v\}.$$

3. To search w^* , compute $\text{LSH}_1(w^*), \dots, \text{LSH}_\beta(w^*)$, and retrieve $\cup_{j=1}^\beta M_j[\text{LSH}_j(w^*)]$.

The maps consists of $y_i, \{w_i\}$ pairs. The keywords y_i are placed into the map, sorted (lexographically) and used as nodes in a binary tree (along with $\{w_i\}$). Internal nodes have store the minimum keyword in the right subtree and the location of the two children LC, RC. We show this design in Figure 6. Let $\text{Tr}_1, \dots, \text{Tr}_\beta$ be the output of BIndex on maps M_1, \dots, M_β respectively.

To turn this into an oblivious search algorithm, one can place each tree in a distinct ORAM. The construction fully traverses every tree Tr_i meaning that there is a constant number of accesses to each ORAM with every search.

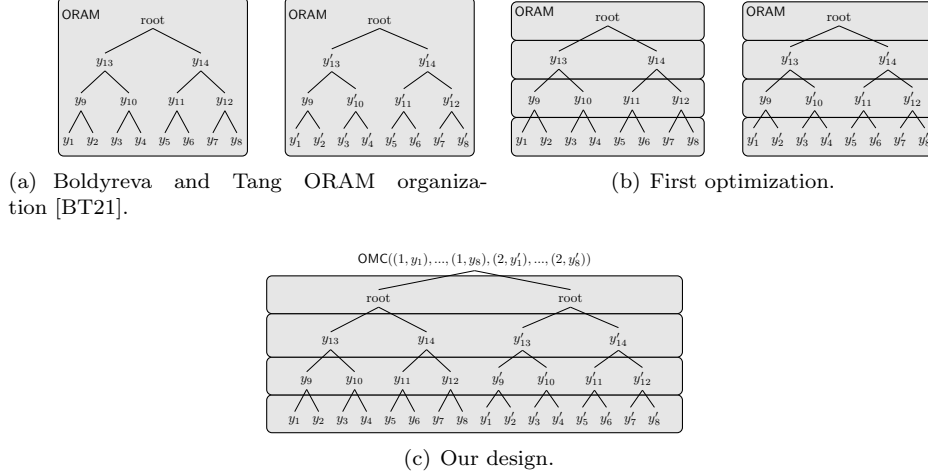


Figure 5: ORAM Organization Strategies. Each shaded region represents data stored together in a single oblivious RAM. In (a) one uses a separate ORAM for each LSH. In (b) we split this to one ORAM per tree level. Then (c) is our design of DOMapE where each level across binary trees is stored in a single ORAM. This allows one to hide which trees are accessed.

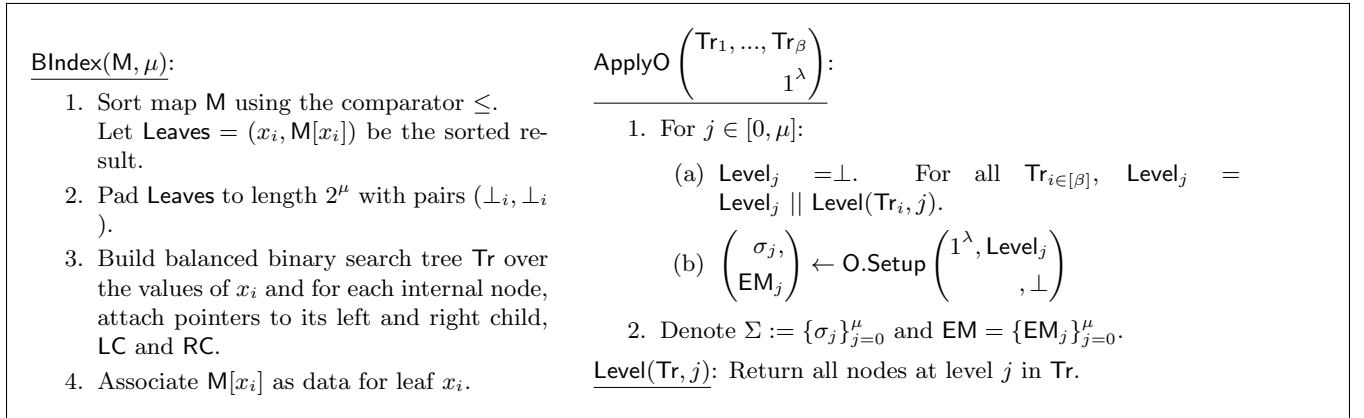


Figure 6: Build tree index and apply ORAM algorithms.

Let μ_i be the number of elements in M_i , define $\mu = \lceil \log \max_i \mu_i \rceil$, by padding each ORAM to length 2^μ each ORAM receives exactly $\mu + 1$ accesses with each query $((\mu + 1)\beta$ across the β trees). This corresponds to Boldyreva and Tang’s approach and is shown in Figure 5(a). However, since each level of each map receives a single access per query, one can store each level of the tree in a separate ORAM. This design is shown visually in Figure 5(b), with each shaded region representing a separate ORAM. With this approach, each shaded region sees exactly one read for each search query. This organization allows higher levels of the tree to serve as the position map of their child eliminating the need for a recursive ORAM [WNL⁺14].

Our approach Recall that our goal is a two part construction: First one queries the OMC to find out which $\delta \leq \beta$ LSHs have matches. Then one queries the relevant δ maps M_i to find records. In this new design, one does not query every M_1, \dots, M_β . As such, the set of queried maps would be leakage. We merge the ORAMs across maps to prevent this. However, we retain a separate ORAM for each level of the trees. This is shown visually in Figure 5(c) and also described by the ApplyO algorithm in Figure 6. This means that each query now makes δ accesses at each ORAM level. There are $\mu + 1$ levels in total resulting in $\delta(\mu + 1)$ ORAM accesses. In this design, parents store the position map of children, enabling non-recursive ORAM.

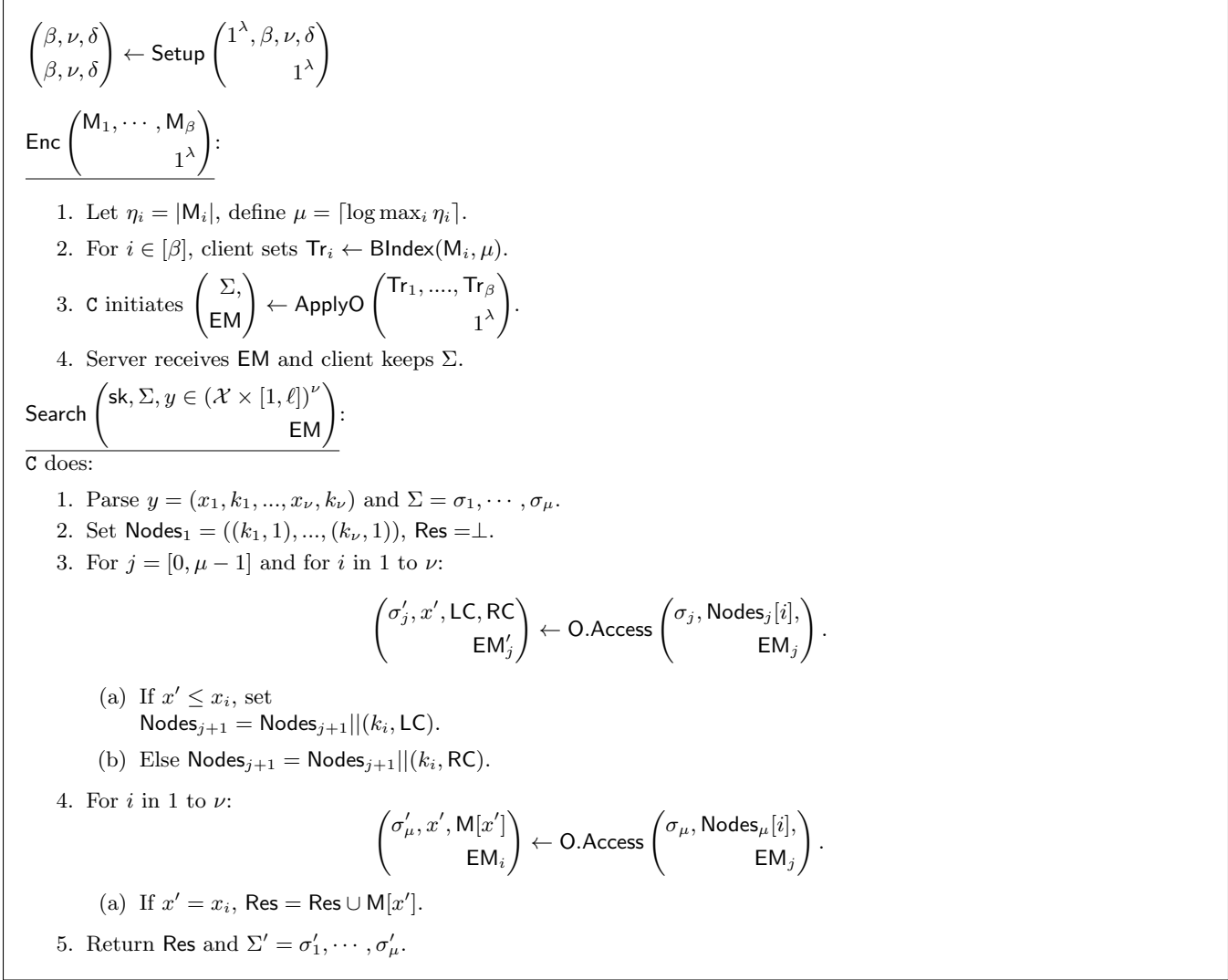


Figure 7: DOMAPE Construction. The BIndex algorithm is shown in Figure 6.

4.2 Detailed design of DOMAPE

Construction 1. Let \mathcal{X} and \mathcal{M} be the domain and range of a map, such that elements in \mathcal{X} are comparable with the \leq operator. Define β maps M_1, \dots, M_β . Let $\text{O} = (\text{O.Setup}, \text{O.Access})$ be an oblivious RAM as defined in Definition 2, and let O_i denote its instantiation for level $0 \leq i \leq \mu$. Consider the DOMAPE construction shown in Figure 7.

Theorem 1. For any $\delta, \beta \in \mathbb{N}$ where $\delta \leq \beta$. Construction 1 describes an DOMAPE for $\mathcal{L}_{\text{BIndex}}(M_1, \dots, M_\beta) = \mu = \max_i \lceil \log |M_i| \rceil$ for $\beta - \text{nmaps}$, $\mu - \text{msize}$, and $\delta - \text{qsize}$.

Proof. We need to show that for every adversary $\mathcal{A}_{\text{DOMAPE}}$ there exists simulator $\text{Sim}_{\text{DOMAPE}}$ for DOMAPE such that $\text{Real}_{\mathcal{A}_{\text{DOMAPE}}}(\lambda) \approx \text{Ideal}_{\mathcal{A}, \text{Sim}_{\text{DOMAPE}}}(\lambda)$ where $\text{Real}_{\mathcal{A}_{\text{DOMAPE}}}$ and $\text{Ideal}_{\mathcal{A}, \text{Sim}_{\text{DOMAPE}}}$ are defined as in Definition 3.

The $\text{Sim}_{\text{DOMAPE}}$ gets as input μ which it uses to initialize that many levels of ORAM simulators. Let Sim_{O_j} denote the simulator for the j^{th} level ORAM, $0 \leq j \leq \mu$. We build $\text{Sim}_{\text{DOMAPE}}$, the simulator for DOMAPE, as follows:

1. Receive inputs $(q, 1^\lambda)$ and $\text{param} = (\beta = \text{nmaps}, \mu = \text{msize}, \delta - \text{qsize})$.

2. For $1 \leq i \leq q \cdot \delta$:
 - (a) For $0 \leq j \leq \mu$, run O_j simulator $(\text{ts}_{O_j,0}^i, \dots, \text{ts}_{O_j,\beta}^i) \leftarrow \text{Sim}_{O_j}(\beta, |\text{Level}_j|, 1^\lambda)$.
 - (b) Set $\text{ts}^{i+1} = (\text{ts}_{O_0,1}^i, \dots, \text{ts}_{O_\mu,\beta}^i)$.
3. Set $\text{ts}^0 = \perp$ and $\text{ts}^1 = (\text{ts}_{O_0,0}^1, \dots, \text{ts}_{O_\mu,0}^1)$.
4. Return $(\text{ts}^0, \dots, \text{ts}^{q \cdot \delta + 1})$.

We then use a hybrid argument where at each step, we replace an ORAM by its corresponding simulator. We obtain the following games

- *Game 0*: O_0, \dots, O_μ ,
- *Game j* : $\text{Sim}_{O_0}, \dots, \text{Sim}_{O_j}, O_{j+1}, \dots, O_\mu$,
- *Game $\mu + 1$* : $\text{Sim}_{O_0}, \dots, \text{Sim}_{O_\mu}$,

with $0 \leq j \leq \mu$.

Note that Game 0 contains $\mu + 1$ ORAM instantiations, which corresponds to the real world $\text{Real}_{\mathcal{A}, \text{DOMapE}}$. Also note that Game $\mu + 1$, contains $\mu + 1$ ORAM simulators, which is equivalent to $\text{Sim}_{\text{DOMapE}}$ and to the ideal world $\text{Ideal}_{\mathcal{A}, \text{SimDOMapE}}$. Then for each Game, we show indistinguishability with the previous one by relying on the security of the underlying O_i .

Lemma 1. *For $1 \leq i \leq \mu + 1$, Game i is indistinguishable from Game $i - 1$.*

Proof. By security of O_i , we have $O_i \approx \text{Sim}_{O_i}$. Since Game $i - 1$ and Game i only differ at index i , we conclude that these two games are indistinguishable. \square

By applying Lemma 1 to each game, we obtain that Games 0 and $\mu + 1$ are indistinguishable, which implies that $\text{Real}_{\mathcal{A}, \text{DOMapE}}$ and $\text{Ideal}_{\mathcal{A}, \text{SimDOMapE}}$ are also indistinguishable and concludes this proof. \square

4.3 OSE design

Construction 2. *For a database $\mathcal{DB} = (w_1, \dots, w_\ell)$ define $\mu = \lceil \log \ell \rceil$. Fix parameters $\delta, \nu, \beta \in \mathbb{N}$ where $\delta \leq \nu \leq \beta$.*

1. Let DOMapE be a disjunctive oblivious map with encryption with $\beta - \text{nmaps}$, $\mu - \text{msize}$, and $\delta - \text{qsize}$,
2. Let OMC be an oblivious membership check with $2^\mu * \beta - \text{ssize}$, $\beta - \text{qsize}$, and $\delta - \text{rsize}$, and
3. Let \mathcal{H} be a family of locality sensitive hashes.

For a database $\mathcal{DB} = (w_1, \dots, w_\ell)$, define $\text{OSE} = (\text{OSE.Setup}, \text{OSE.Enc}, \text{OSE.Find})$ as in Figure 8.

Theorem 2. *Let DOMapE and OMC be as in Construction 2. Then Construction 2 is an oblivious searchable encryption scheme with leakage $\mathcal{L}_{\text{Enc}}(\mathcal{M}_1, \dots, \mathcal{M}_\beta) = \mu$.*

Proof. Let Sim_{OSE} be the simulator for the OSE scheme. This simulator will run as follows:

1. Upon input $(q, \beta, \delta, \mu, 1^\lambda)$, run the simulator for the oblivious membership check with $\text{param} = (\beta * 2^\mu - \text{ssize}, \beta - \text{qsize}, \delta - \text{rsize}, (\text{ts}_{\text{OMC}}^1, \dots, \text{ts}_{\text{OMC}}^{q+1}) \leftarrow \text{Sim}_{\text{OMC}}(q, \gamma = \beta 2^\mu, \beta, \delta)$.
2. Run $(\text{ts}_{\text{DOMapE}}^0, \dots, \text{ts}_{\text{DOMapE}}^{q+1}) \leftarrow \text{Sim}_{\text{DOMapE}}(q, \text{param}, 1^\lambda)$ with $\text{param} = (\beta - \text{nmaps}, \mu - \text{msize}, \delta - \text{qsize})$.
3. Set $\text{ts}_0^* = \text{ts}_{\text{DOMapE}}^0$.
4. For $1 \leq i \leq q + 1$, set $\text{ts}_i^* = (\text{ts}_{\text{OMC}}^i, \text{ts}_{\text{DOMapE}}^i)$.
5. Output transcripts $\text{ts}_0^*, \dots, \text{ts}_{q+1}^*$.

We use a hybrid argument to show security of OSE. Consider the following games:

1. Real_{OSE} ,

OSE.Setup $\left(\begin{array}{c} 1^\lambda \\ 1^\lambda \end{array} \right)$:

1. Client samples $\text{LSH}_1, \dots, \text{LSH}_\beta \leftarrow \mathcal{H}(1^\lambda)$.
2. Client runs

$$\begin{aligned} \left(\begin{array}{c} \text{sk}_{\text{OMC}} \\ \perp \end{array} \right) &\leftarrow \text{OMC.Setup} \left(\begin{array}{c} 1^\lambda \\ \perp \end{array} \right), \\ \left(\begin{array}{c} \text{sk}_{\text{DMapE}} \\ \text{pp}_{\text{DMapE}} \end{array} \right) &\leftarrow \text{DMapE.Setup} \left(\begin{array}{c} 1^\lambda \\ 1^\lambda \end{array} \right). \end{aligned}$$

3. Denote $\text{pp} = \text{pp}_{\text{DMapE}}$, $\text{sk} = (\text{sk}_{\text{OMC}}, \text{sk}_{\text{DMapE}}, \text{LSH}_1, \dots, \text{LSH}_\beta)$.
4. Client sends pp to server and keeps sk .

OSE.Enc $\left(\begin{array}{c} \text{sk}, \mathcal{DB} = (w_1, \dots, w_\ell) \\ \text{pp}_{\text{DMapE}} \end{array} \right)$:

1. For $1 \leq i \leq \beta$, client:
 - (a) Initializes map M_i .
 - (b) For $1 \leq j \leq \ell$ sets $M_i[\text{keyword}_{i,j}] = \{w_j \mid \text{LSH}_i(w_j) = \text{keyword}_{i,j}\}$.
 - (c) Adds dummy values to M_i until it is of size ℓ .
- 2.

$$\begin{aligned} \left(\begin{array}{c} \perp \\ ES \end{array} \right) &\leftarrow \text{OMC.Enc} \left(\begin{array}{c} \text{sk}_{\text{OMC}}, \cup_{i=1}^{\beta} \cup_{j=1}^{\ell} i \parallel \text{keyword}_{i,j} \\ 1^\lambda \end{array} \right), \\ \left(\begin{array}{c} \sigma \\ EM \end{array} \right) &\leftarrow \text{DMapE.Enc} \left(\begin{array}{c} \text{sk}, M_1, \dots, M_\beta \\ \text{pp}_{\text{DMapE}} \end{array} \right). \end{aligned}$$

3. Denote $I_c = \sigma$ and $I_s = (ES, EM)$.

OSE.Find $\left(\begin{array}{c} \text{sk}, y, I_c \\ I_s \end{array} \right)$:

1. Let $I_s = (ES, EM)$.
2. Client creates OMC set $EC = (1 \parallel \text{LSH}_1(y), \dots, \beta \parallel \text{LSH}_\beta(y))$.
- 3.

$$\begin{aligned} \left(\begin{array}{c} \text{Res}_{\text{OMC}} \\ \perp \end{array} \right) &\leftarrow \text{OMC.Search} \left(\begin{array}{c} \text{sk}_{\text{OMC}}, EC \\ ES \end{array} \right), \\ \left(\begin{array}{c} r, \sigma' \\ EM' \end{array} \right) &\leftarrow \text{DMapE.Search} \left(\begin{array}{c} \text{sk}_{\text{DMapE}}, \sigma, \text{Res}_{\text{OMC}} \\ EM \end{array} \right). \end{aligned}$$

4. Denote $J = (\cup_{i=1}^{\beta} r_i) \setminus \perp$, $I'_c = \sigma'$ and $I'_s = (ES, EM')$.
5. Client receives J, I'_c and server stores I'_s .

Figure 8: OSE construction from OMC, DMapE and LSH.

2. $\text{Sim}_{\text{OSE}}^*$, which runs Sim_{OMC} and $\text{Real}_{\text{DMapE}}$,
3. Sim_{OSE} , which runs Sim_{OMC} and $\text{Sim}_{\text{DMapE}}$ as described above.

We want to show that game 1 is indistinguishable from game 3.

Lemma 2. *Game 1 and Game 2 are indistinguishable.*

Proof. The difference between games 1 and 2 is $\text{Sim}_{\text{OSE}}^*$'s use of Sim_{OMC} instead of Real_{OMC} . By security of the OMC scheme, game 1 and 2 are indistinguishable. \square

Lemma 3. *Game 2 and Game 3 are indistinguishable.*

Proof. The difference between games 2 and 3 is Sim_{OSE} 's use of $\text{Sim}_{\text{DOMapE}}$ instead of $\text{Real}_{\text{DOMapE}}$. By security of the DOMapE scheme, game 2 and 3 are indistinguishable. \square

Combining Lemmas 2 and 3, we obtain that games 1 and 3 are indistinguishable, which conclude our proof. \square

Theorem 2 does not handle correctness. Since there is an overlap between the histograms for real data in Figure 11 one cannot make strong correctness claims. We evaluate correctness empirically in Section 7.

4.4 Oblivious membership check constructions

We discuss options to implement OMC. We briefly cover approaches based on Bloom filter lookups. In Appendix 4.4.1, we describe how to build OMC from private set intersection. This is the tool that we use for microbenchmarks. In our implementation, we use a local Bloom filter to emulate an OMC.

Oblivious Bloom Filter Lookups The client's set can be stored in a Bloom [Blo70], Cuckoo [FAKM14], or XOR [GL20] filters which is then stored on the server in an ORAM. The client will request the relevant bits from the ORAM. This prevents the client from having to store the entire filter on their side, but requires them to request multiple ORAM accesses to query the relevant bits. BlindSEER [PKV⁺14, FVK⁺15] built a tree of encrypted Bloom filters for general Boolean search. Search of each node uses Garbled circuits to decide whether to proceed to children. One can use a single level of their tree as an OMC as long as only the client learns the response. This requires some modification as their system was optimized for circuits that output a bit, we would need the set of matching locations. Their system was evaluated on datasets with 10^8 records [FMC⁺15].

4.4.1 Building OMC from PSI and pseudorandom permutations

Private set intersection (PSI) [FNP04] is a form of secure multi-party computation where a client and server hold sets Y and X respectively. They run an interactive computation, at the end, the client learns $X \cap Y$. No other information is leaked. Current implementations of PSI depend on one of two tools: oblivious polynomial evaluations (OPE) and oblivious pseudorandom functions (OPRF).⁷ We show how to build OMC from honest-but-curious PSI as follows:

1. At initialization the client applies a pseudorandom permutation (PRP) to each element in the set X .
2. The client sends the set of elements (passed through the PRP) to the server.
3. Later when the client has a set Y , they apply the pseudorandom permutation to each element of Y , and uses the resulting values as their set for the PSI protocol.

In OMC, the simulator learns the size of both sets X, Y , using an ideal PSI, only the size of X is leaked to the server. Both the sizes of X and Y are global parameter, $\beta \cdot 2^\mu$ and β respectively. We now formalize the above, re-stating the PSI definition used for VOLE-PSI [RS21, Figure 5].

Definition 7 (Private Set Intersection (PSI)). *Let \mathcal{X} denote a set. Let the client hold a set $X \subset \mathcal{X}$ and the server hold a set $Y \subseteq \mathcal{X}$. Consider $\begin{pmatrix} Z \\ \perp \end{pmatrix} \leftarrow \text{PSI} \begin{pmatrix} X \\ Y \end{pmatrix}$, the PSI protocol between the client and the server. The client learns Z and the server learns nothing.*

⁷Cristofaro et al. [CT10] construct efficient PSI protocols, under the restriction that the server can do some precomputation or the client is weak. Dong et al. [DCW13] present an efficient PSI protocol based on a variant of Bloom Filters called *garbled Bloom Filters*. Dachman-Soled et al. [DSMRY09] present an efficient PSI protocol utilizing secret sharing and Reed-Soloman codes. Malicious secure implementations of PSI utilizing OPE also depend on zero-knowledge proofs to prevent parties from deviating from the protocol.

<p><u>Real$_{\mathcal{A},n_X,n_Y}(1^\lambda)$:</u></p> <ol style="list-style-type: none"> 1. $(X, Y) \leftarrow \mathcal{A}(1^\lambda)$ where $X, Y \subset \mathcal{X}$. 2. If $X > n_X$ or $Y > n_Y$, return \perp. 3. Output ts of $\begin{pmatrix} Z, \\ \perp \end{pmatrix} \leftarrow \text{PSI} \begin{pmatrix} X, \\ Y \end{pmatrix}.$	<p><u>Ideal$_{\mathcal{A},\text{Sim},n_X,n_Y}(1^\lambda)$:</u></p> <ol style="list-style-type: none"> 1. $(X, Y) \leftarrow \mathcal{A}(1^\lambda)$ where $X, Y \subset \mathcal{X}$. 2. Output $\text{ts} \leftarrow \text{Sim}(1^\lambda, \text{param})$ where $\text{param} = (X = \text{csize}, Y = \text{ssize})$.
---	--

Figure 9: Definition of PSI security.

<p><u>OMC.Setup</u> $\begin{pmatrix} 1^\lambda \\ \perp \end{pmatrix}$:</p> <ol style="list-style-type: none"> 1. Client randomly samples a PRP key $\text{sk} \xleftarrow{\\$} \{0, 1\}^\lambda$. <p><u>OMC.Enc</u> $\begin{pmatrix} \text{sk}, X \\ \perp \end{pmatrix}$:</p> <ol style="list-style-type: none"> 1. Client initializes $ES = \emptyset$. 2. For each $x_i \in X$, client sets $ES = ES \cup F(\text{sk}, x_i)$, where F is a PRP. 3. Client sends ES to server. 	<p><u>OMC.Search</u> $\begin{pmatrix} \text{sk}, Y \\ ES \end{pmatrix}$:</p> <ol style="list-style-type: none"> 1. For each $y_j \in Y$, client computes permutation $F(\text{sk}, y_j)$. 2. Client initializes a map M and inserts $M[F(\text{sk}, y_j)] = y_j$. 3. Client sets $EC = \{F(\text{sk}, y_j)\}$. 4. Run $\begin{pmatrix} EC \cap ES, \\ \perp \end{pmatrix} \leftarrow \text{PSI} \begin{pmatrix} EC, \\ ES \end{pmatrix}$. 5. Client <ol style="list-style-type: none"> (a) Initializes $\text{Res} = \emptyset$. For each $c \in EC \cap ES$, set $\text{Res} = \text{Res} \cup M[c]$. (b) If $\text{Res} > \delta$ set Res to be δ random elements of Res. (c) If $\text{Res} < \delta$ append elements from \perp_1, \dots to make Res size δ. (d) Outputs Res.
---	--

Figure 10: Construction of OMC from PSI and PRP.

Correctness: *Correctness is,* $\Pr \left[Z \neq X \cap Y \mid \begin{pmatrix} Z, \\ \perp \end{pmatrix} \leftarrow \text{PSI} \begin{pmatrix} X, \\ Y \end{pmatrix} \right] \leq \text{negl}(\lambda)$.

Security: *PSI is secure if for any PPT adversary \mathcal{A} , there exists a simulator Sim , such that the distributions $\text{Real}_{\mathcal{A},q}$ and $\text{Ideal}_{\mathcal{A},\text{Sim},q}$, described in Figure 9, are computationally indistinguishable. VOLE-PSI [RS21] uses a different security definition for PSI, their definition implies ours.*

Remark *We consider PSI where both X and Y are constant size sets so we ignore the case when a party provides too large a set.*

Definition 8 (Pseudorandom Permutation (PRP)). *Let $F : \{0, 1\}^\lambda \times \mathcal{X} \rightarrow \mathcal{X}$ be an efficient keyed permutation and \mathcal{F}_n denote the set of all permutations on \mathcal{X} . F is a pseudorandom permutation if for any PPT adversary \mathcal{A} ,*

$$\left| \Pr \left[\mathcal{A}^{F_k(\cdot)}(1^\lambda) = 1 \mid k \xleftarrow{\$} \{0, 1\}^\lambda \right] - \Pr \left[\mathcal{A}^{f(\cdot)}(1^\lambda) = 1 \mid f \xleftarrow{\$} \mathcal{F}_n \right] \right| \leq \text{negl}(\lambda).$$

Construction 3 (OMC from PSI and PRP). *Let PSI denote a PSI scheme and $F : \{0, 1\}^\lambda \times \mathcal{X} \rightarrow \mathcal{X}$ be a pseudorandom permutation. Let X and Y be sets, such that $X \subseteq \mathcal{X}$ and $Y \subseteq \mathcal{X}$. Then we build OMC as in Figure 10.*

Theorem 3. *Let PSI be a secure private set intersection scheme and F be a pseudorandom permutation, then Construction 3 describes a secure OMC with $\gamma - \text{ssize}$, $\beta - \text{qsize}$, and $\delta - \text{rsize}$.*

Proof. Correctness is straightforward and follows from correctness of the PSI scheme and the fact that each element in F is a permutation. If there are more than δ elements in the intersection then δ are chosen randomly. However,

the correctness guarantee only requires that non- \perp elements returned be in the intersection. Security follows from the security of the PRP (values seen by the server are indistinguishable from random) and the security of the PSI scheme (server learns nothing about EC and $EC \cap ES$). Formally, we build the simulator Sim_{OMC} as follows:

1. Receive inputs 1^λ , number of queries q , server’s set size γ .
2. For $1 \leq i \leq \gamma$, sample a random value $x_i \xleftarrow{\$} \mathcal{X}$.
3. Define the server set as $ES = \{x_i\}_{i=1}^\gamma$.
4. For $1 \leq k \leq q$: run the PSI simulator $\text{ts}^k \leftarrow \text{Sim}_{\text{PSI}}(1^\lambda, \text{param})$ with $\text{param} = (\beta = \text{csize}, \gamma = \text{ssize})$.
5. Output ES and $\text{ts}^1, \dots, \text{ts}^q$.

We then use a hybrid argument to show security of OMC. Consider the following games:

1. Real_{OMC} ,
2. Replace the PRP F by a random permutation $f \xleftarrow{\$} \mathcal{F}_n$,
3. Sim_{OMC} .

Games 1 and 2 are indistinguishable by the security of the PRP scheme, and games 2 and 3 are indistinguishable by the security of the PSI scheme. We note that sampling a random output value is equivalent to sampling a random input value as inputs are guaranteed to be a set. \square

5 Datasets

We test and evaluate our implementation on three datasets:

ND-0405 dataset This dataset [PSO⁺09,BF16] is a superset of the NIST Iris Evaluation Challenge [PBF⁺08]. It consists of the readings of left and right irises from 356 individuals, each iris having at least 4 distinct readings. We use the ThirdEye feature extractor [AF19] to obtain 1024 bits feature vectors from the original iris images. Since the left irises were used to train the feature extractor, we use the right ones for testing and evaluation. The first reading of each right iris is in the \mathcal{DB} ; queries come from the remaining readings. The Hamming distance distributions are in Figure 11(a).

Synthetic dataset Available irises datasets are of limited size, often no more than a few hundreds irises (356 individuals for ND-0405). Real world systems would store thousands to millions individuals, depending on the application. Our solution is to generate *synthetic* irises templates, that mimic actual ones. As can be seen in Figure 11b, synthetic data same and different distributions are similar to the ND-0405 ones. The details on synthetic data generation are in Appendix A. The high level approach is a generative adversarial network (GAN) [GPAM⁺20] as in prior approaches on synthetic iris generation [AF20].

Random dataset This dataset is made from randomly generated 1024 bits vectors. The Hamming distance between two vectors is close to 0.5 with a small variance. This is visible in the red histogram from Figure 11c.

Random and synthetic queries generation Contrary to the ND0405 dataset [BF16], the random and synthetic datasets do not include queries. We generate queries from a distribution that resembles the one for ND-0405. We use the common observation that like irises comparisons have a distribution close to a binomial across different feature extractors [Dau09,Dau05,SSF19]. From Figure 11a, we extract the mean, $\mu = 0.21$, and the standard deviation, $\sigma = 0.056$. This yields a distribution $B(n, \mu)/n$, the binomial distribution for $n = 53$. This is because for $B(n, \mu)$ it is true that $\sigma^2 = \mu(1 - \mu)n$. Thus, by linearity of expectation for $B(n, \mu)/n$ it is true that $\sigma^2 = \mu(1 - \mu)/n$, thus one can compute $n = \lceil \mu(1 - \mu)/\sigma^2 \rceil = \lceil 52.9 \rceil = 53$.

We then generate queries for the random and synthetic datasets as follows:

1. Generate a binomial distribution using the mean and standard deviation of the same iris distribution for the ND-0405.

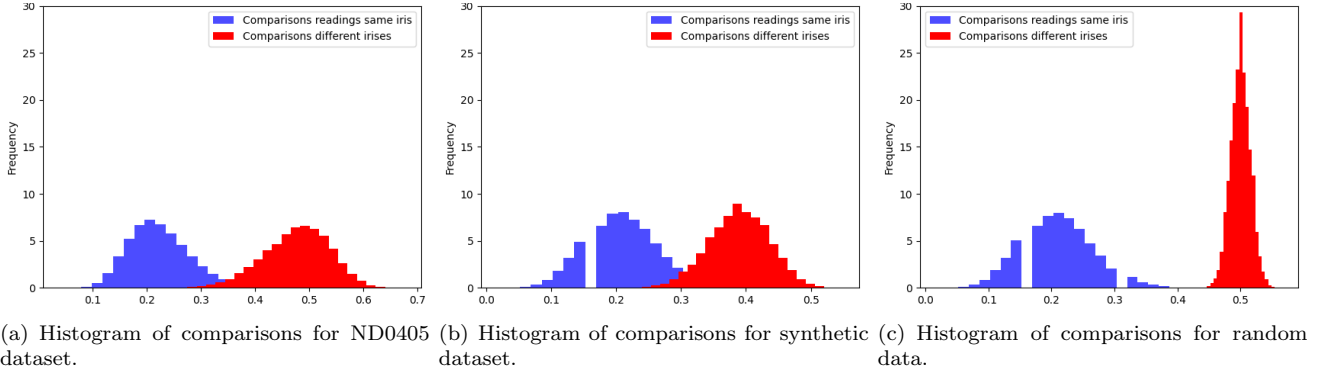


Figure 11: Histograms of Hamming distance between readings of the same iris (in blue) and different irises (in red). Different irises are stored in the database and queries are drawn from a different reading of an iris in the database. The gaps in the synthetic and random blue histograms are caused by the query generation technique used (see paragraph on random and synthetic queries generation) in Section 5.

2. For each feature vector in the dataset, create a corresponding query by sampling an error fraction from the $\text{frac} \leftarrow B(53, 0.21)/53$.
3. Flip the number of error bits, $\text{frac} * 1024$, in the feature vector.

Using this technique, we obtain the same iris distributions (in blue) for synthetic and random data shown in Figures 11b and 11c. There are only 54 possible outcomes for a fraction of error bits, this leads to discontinuities in the histograms presented in Figures 11b and 11c.

6 Implementation

We present an open-source implementation of our algorithms including the LSH parameter finding, tree building, and oblivious search [pri24]. This implementation is in Python 3.10 and uses the PathORAM [SDS⁺18] module [Hac18]. Our experiments use a Bloom filter cache on the client as an OMC to focus on the performance of the developed DOMapE. We separately evaluate an OMC candidate based on PSI in Section 7. Our implementation supports two main conclusions.

1. One can set a $\delta < \beta$ size of the query to DOMapE that supports a high true accept rate. For a query, we define *bad matches* to be the number of LSH matches that only result in incorrectly returned records; setting δ to be 1 more than the 95% of this value. See Table 2 for a comparison of true accept rate for the setting when $\delta = \beta$ and when $\delta < \beta$. In all analyzed parameters $\delta/\beta < .06$. The value δ is higher for real and synthetic data than for random data; this is due to the larger variance of distances between readings of different irises. This increases the number of false matches. As one exception, we set δ for our 25K random parameters heuristically based on the smaller dataset sizes.
2. The two stage DOMapE approach improves search performance. While setup takes several hours, parallel search is at most 35ms. Without the two stage approach parallel search would require a server with hundreds of threads. See further timing discussion in Section 7.2.

We see two mechanisms for achieving better performance as datasets grow:

1. Use ORAM that supports sending multiple queries in parallel (this is a weaker object than a parallel ORAM [WST12, BCP16, CLT16]).
2. Group LSHs into groups that are placed into a single RAM and analyze the required δ for each group. As $\delta/\beta \approx .06$ one can use standard concentration bounds to argue about δ for each group as long as the data and queries are independent of the LSHs.

Dataset Modifications Our implementation does not allow for insertion after the initial building of the tree. With ORAM one can rebuild the trees using techniques of Wang et al. [WNL⁺14].

Alternative to ORAM In the static setting, one can use private information retrieval (PIR) [CG97, CKGS98] with encryption. At retrieval, single server computational PIR and PathORAM [SDS⁺18] with “large” blocks of size $\Omega(\log^2 N)$ both achieve communication complexity of $O(\log N)$, with N the number of blocks. However, time efficiency would probably suffer from this change. Traditional PIR schemes require work $\Theta(|\mathcal{DB}|)$ on the server. Doubly efficient PIR (DEPIR) [BIM00] preserves the communication efficiency of regular PIR but with $o(|\mathcal{DB}|)$ server work. To achieve this DEPIR relies on a server pre-processing stage which is allowed in our model. DEPIR constructions are based on ring LWE [LMW23] or a non-standard secretly permuted Reed-Muller codes assumption [BIPW17, CHR17]. Currently, the ring LWE DEPIR is asymptotically efficient but inefficient in practice [OPPW24].

7 Evaluation

Evaluation is split into two parts: 1) parameter analysis and accuracy, and 2) efficiency of the resulting cryptographic construction. Our parameter analysis focuses on the TAR and number of matches. Our efficiency analysis focuses on network roundtrips, storage, and single-threaded computation time.

7.1 Accuracy - Parameter analysis

Each experiment is conducted on each dataset. Recall the relevant parameters: α , the length of the extended LSH, $\beta - \text{mmaps}$, and $\delta - \text{qsize}$.

Finding parameters The first part of the experiment was a manual search across parameters α, β , measuring the TAR and number of bad matches. Selected parameters had TAR of at least 90%. The average number of bad matches was at most 10 for random data and at most 50 for ND and synthetic data. Once α, β were selected we recorded the histogram of bad matches and set δ to be one more than the 95% of this histogram (Fig. 12).

Measuring accuracy We then measured accuracy for a search that queries all β maps and one that only queries δ maps. For these tests, we only measure the TAR to understand the impact of restricting the number of searched values on accuracy.

As one exception, due to slow speed and high memory overhead for the 25K random dataset, we picked $\delta = 20$ based on our parameters from smaller datasets, this dataset naturally had a larger number of bad matches but still displayed a high TAR of .91 when restricted to $\delta = 20$ maps.

Discussion In proximity search, high TAR requires capturing the tail of comparisons between different readings of the same iris (shown in Figure 11). For example, for distance $t = .21n$ and a FAR of .01, Section 3.1 proposed $\beta = 65$ and $\alpha = 13$. Table 2 shows that even for random data, we require $\beta = 630$ and $\alpha = 15$. These parameters increase further on the ND and Synthetic datasets. Parameters vary between random and synthetic data due to a larger variance in the Hamming distance histograms, shown in Figure 11. This is in contrast to the histogram of good matches across datasets in Figure 13 which are consistent across datasets. This leads to an increase in the selected δ . Across dataset sizes, δ for synthetic data is about 5 times δ for random data. The ND and synthetic data statistics align well. This gives some indication that parameters for larger synthetic dataset sizes would yield comparable performance on real irises. Across all parameter settings $\delta/\beta < .06$ validating the overall design. Restricting to only δ accesses in the DOMapE does harm TAR. Note that the real system does not know which matches are good so even in the case when many matches are bad it is possible for a single good match to be included in the arbitrarily selected δ traversals. The worst degradation is for synthetic data with 5000 records where TAR drops from .91 to .83.

7.2 Speed - Cryptographic Efficiency

The implementation was tested on a AMD Ryzen Threadripper PRO 7995WX CPU with 96 cores and 768 GB of RAM, running Ubuntu 22.04. Results are in Table 3. We did not model network delay.

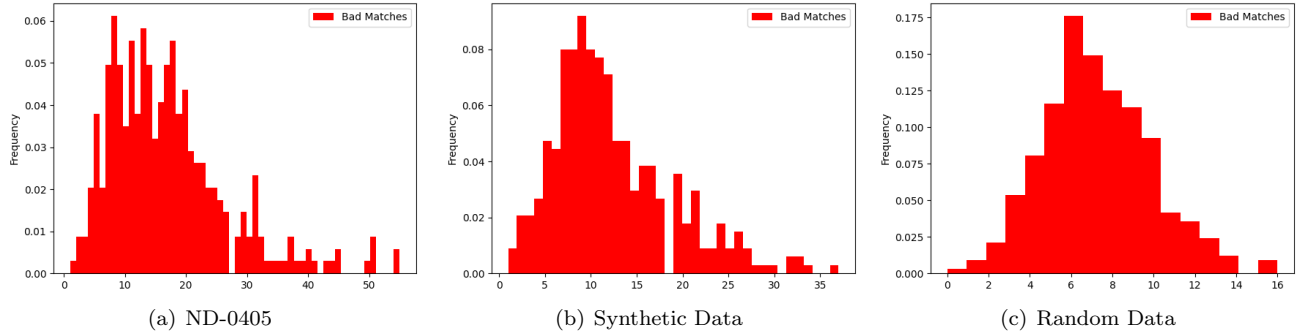


Figure 12: Number of bad matches across datasets of size 356.

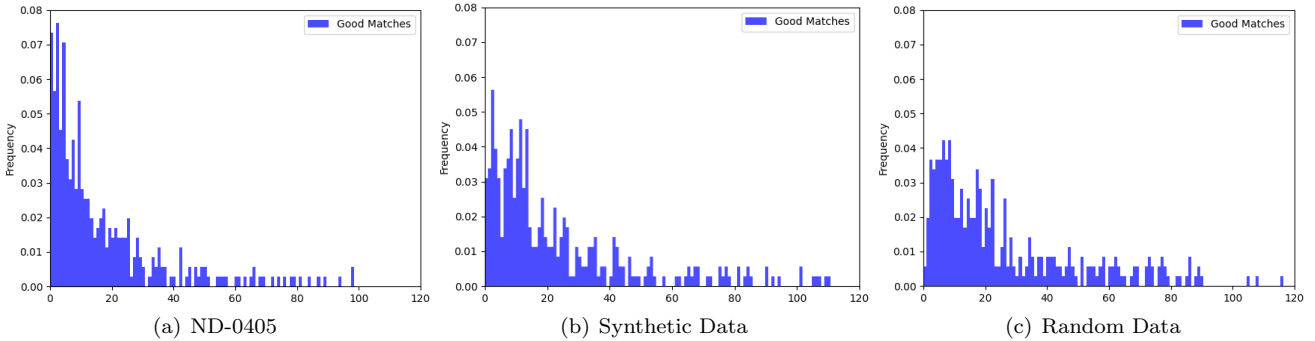


Figure 13: Number of good matches across datasets of size 356.

Storage efficiency Feature vectors are 1024 bit vectors, so 5K irises is 640 KB. The unprotected (same structure as DOMapE but without ORAM) index takes approximately 122.3 MB.⁸ This represents a storage increase factor of around 22 between raw data and unprotected index. As shown in Table 3, for our encrypted storage this amounted to 35.6 GB in storage. ORAM increases storage again approximately 291 times. As we discuss in the Conclusion, one can more efficiently pack ORAM blocks using trees with a branching factor > 2 .

Time efficiency The time to build the encrypted index is largely dominated by the ORAM setup time so we only report the later (column “O.Init”). For small datasets (356 records) ORAM setup takes hours, while larger datasets take days.

Parallel search time, the sum of the max required time to complete each ORAM read at each level, remains under 35ms on all tested parameters. We note that parallel search time is less than sequential search time divided by δ as sequential search time includes reading all ORAMs from disk at the beginning of each search and writing them back to disk at the end of each search. This is done with every query. Only searching δ trees is critical to process the ORAM reads in parallel on a moderately powerful modern server.

Network Round Trips We report on two figures, the number of round trips using a purely sequential non-recursive PathORAM implementation and the number of roundtrips if one is able to fully batch all requests at the same level. One can store the position map for ORAMs in the prior tree level [WNL⁺14], enabling non-recursive constructions. For the largest synthetic parameter sizes, if one assumes a fast network with 60ms responses and unbounded bandwidth then network delays result in .84 seconds in parallel rounds trips, but slower 1s responses results in 14 seconds.

⁸Internal node consists of an LSH number, a 22 LSH values, and 2 child identifiers (either the node id or the position of the child in the next ORAM). Leaf nodes consist of a single 32 node identifier.

Dataset size	Dataset type	α	β		Matches			TAR		
			nmaps	# FA	Avg bad	Max bad	Avg good	without OMC	δ qsize	with OMC
356	random	15	630	7.5	7.2	16	32.1	0.98	13	0.94
356	ND	18	850	14.4	16.8	55	21.6	0.95	37	0.89
356	synthetic	18	850	15.5	12.2	37	22.6	0.96	26	0.92
1000	random	18	850	3.6	3.5	11	21.2	0.96	8	0.96
1000	synthetic	19	1000	3.7	22.7	82	20.2	0.95	47	0.96
2500	random	19	1000	5.7	5.6	13	25	0.94	11	0.89
2500	synthetic	21	1200	16	38.4	569	6.6	0.92	56	0.85
5000	random	20	1200	6.9	6.7	14	21.7	0.92	12	0.87
5000	synthetic	22	1300	21.4	44.7	578	6.9	0.91	72	0.83
25000	random	22	3500	8.7	57.2	277	6.1	-	20	0.91

Table 2: TAR/FAR and the number of matches for random, ND0405, and synthetic datasets of different sizes. For 25K dataset, δ was set manually (not 95% of bad matches).

Dataset size, ℓ	Dataset type	β		δ	# Queries	ORAM Reads	# Roundtrips		Time (s)			Size EDB (GB)
		α	nmaps	qsize			seq.	par.	O.Init	seq.	par.	
356	random	15	630	13	356	117	118	10	1.1×10^3	.65	.013	1.1
356	ND	18	850	37	356	333	334	10	1.5×10^3	.91	.013	1.1
356	synthetic	18	850	26	356	234	235	10	1.6×10^3	.85	.013	1.1
1000	random	18	850	8	500	80	81	11	3.1×10^3	1.06	.015	2.2
1000	synthetic	19	1000	47	500	470	471	11	3.6×10^3	1.51	.016	2.2
2500	random	19	1000	11	500	132	133	13	16×10^3	3.44	.018	8.9
2500	synthetic	21	1200	56	500	672	673	13	20×10^3	4.52	.020	17.8
5000	random	20	1200	12	500	156	157	14	42×10^3	7.41	.021	35.6
5000	synthetic	22	1300	72	500	936	937	14	45×10^3	8.82	.023	35.6
25000	random	22	3500	20	100	300	301	16	552×10^3	7.89	.035	285

Table 3: Efficiency results. O.Init is time to initialize all ORAMs. O.Read is average read time (across ORAM layers). Search is time per query and includes tree traversals. Size EDB denotes the size of the ORAM files that are stored on the server (OMC storage is ignored since it is much smaller). Sequential number of roundtrips is $1 + \# \text{ORAM Reads}$ and Parallel Rounds trips is $\lceil \log_2 \ell \rceil + 1$. All timing numbers are averaged across the number of queries in # Queries.

Using a batched ORAM implementation Many existing ORAM schemes including PathORAM naturally supported batched read/write operations where the client keeps a larger stash. In the case of PathORAM, the client repeatedly reads and writes a “random” path on a tree. One can naturally perform all reads first and then perform all writes, simulating the intermediate storage that would be held by the server. Parallel ORAM is a more complex solution when the reads come from different clients [WST12].

Evaluation of OMC implementation using private set intersection On the same hardware as the rest of the evaluation we deployed the VolePSI implementation [RS21]. To test the largest synthetic parameters, we deployed this with a server set of size 6.5 million items and a client set of 1300 items. This corresponds to the largest set of parameters in Table 3. VolePSI is based on OT extension and requires a setup phase. We benchmarked 32 PSI iterations with the first taking 766ms and the rest taking 2ms of computation. We note that VolePSI requires 7 messages of communication. To get the results in our abstract and introduction, we add 2ms to Table 3 and four rounds of communication. These results justify the focus on the design of DOMapE.

8 Conclusion

Private Eyes was tested with parallel response times of at most 35ms on databases of thousands of irises. Our construction combines LSHs and oblivious maps. The unique aspect of our design is the recognition and mitigation of the cryptographic inefficiencies caused by the high noise in biometric data. The statistics of biometric data inspired a two-stage approach which filters which LSHs to query using a lighter-weight membership checking primitive before the heavy-weight oblivious map.

We used binary trees but one could use trees with a higher branching factor (or skiplists as in [BC14]) to reduce the number of ORAM lookups. Ideally, each node would correspond to a single ORAM block which is commonly a multiple of 256 bytes. Our current estimate is that internal nodes account for ≤ 128 bits of storage out of the 256 byte block size. As such, one could make the tree into a 18-ary tree (32 bits for LSH number, 32 bits for left most child, and $22 + 32$ bits for each additional comparison node). This would reduce the depth of the trees and required number of round trips by a factor of $\log_2 18 \approx 4$.

Acknowledgements

The authors thank the anonymous reviewers for their help in improving the manuscript. The authors also thank Mayank Varia for helpful discussions and Mike Roseluk for helping provide options for the evaluation of private set intersection protocols. J.H. is supported by NSF Award #1950600. C.C. is supported by NSF Award #2141033. L.D. is supported by the Harriott Fellowship. S.A. is supported by a fellowship from Synchrony Inc. and the State of Connecticut. B.F. is supported by NSF Awards # 2141033 and 2232813 and ONR Grant N00014-19-1-2327.

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via Contract No. 2019-19020700008. This material is based upon work supported by the Defense Advanced Research Projects Agency, DARPA, under Air Force Contract No. FA8702-15-D-0001. Any opinions, findings, conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of DARPA. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government.

The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

References

- [ACD⁺22] Sohaib Ahmad, Chloe Cachet, Luke Demarest, Benjamin Fuller, and Ariel Hamlin. Proximity searchable encryption for the iris biometric. In *AsiaCCS*, 2022. <https://ia.cr/2020/1174>.
- [AF19] Sohaib Ahmad and Benjamin Fuller. Thirdeye: Triplet-based iris recognition without normalization. In *IEEE BTAS*, 2019.
- [AF20] Sohaib Ahmad and Benjamin Fuller. Resist: Reconstruction of irises from templates. In *IJCB*, pages 1–10. IEEE, 2020.
- [BBC⁺10] Mauro Barni, Tiziano Bianchi, Dario Catalano, Mario Di Raimondo, Ruggero Donida Labati, Pierluigi Failla, Dario Fiore, Riccardo Lazzeretti, Vincenzo Piuri, Fabio Scotti, and Alessandro Piva. Privacy-preserving fingercode authentication. In *Proceedings of the 12th ACM Workshop on Multimedia and Security*, page 231–240, New York, NY, USA, 2010. Association for Computing Machinery.
- [BBOH96] Christopher M Brislawn, Jonathan N Bradley, Remigius J Onyszczyk, and Tom Hopper. The FBI compression standard for digitized fingerprint images. In *Proc. SPIE*, volume 2847, pages 344–355, 1996.
- [BC14] Alexandra Boldyreva and Nathan Chenette. Efficient fuzzy search on encrypted data. In *International Workshop on Fast Software Encryption*, pages 613–633. Springer, 2014.
- [BCP16] Elette Boyle, Kai-Min Chung, and Rafael Pass. Oblivious parallel ram and applications. In *Theory of Cryptography Conference*, pages 175–204. Springer, 2016.

- [BF16] Kevin W Bowyer and Patrick J Flynn. The ND-IRIS-0405 iris image dataset. *arXiv preprint arXiv:1606.04853*, 2016.
- [BG11] Marina Blanton and Paolo Gasti. Secure and efficient protocols for iris and fingerprint identification. In Vijay Atluri and Claudia Diaz, editors, *Computer Security – ESORICS 2011*, pages 190–209, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [BHJP14] Christoph Bösch, Pieter Hartel, Willem Jonker, and Andreas Peter. A survey of provably secure searchable encryption. *ACM Computing Surveys (CSUR)*, 47(2):1–51, 2014.
- [BIM00] Amos Beimel, Yuval Ishai, and Tal Malkin. Reducing the servers computation in private information retrieval: Pir with preprocessing. In Mihir Bellare, editor, *Advances in Cryptology — CRYPTO 2000*, pages 55–73, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.
- [BIPW17] Elette Boyle, Yuval Ishai, Rafael Pass, and Mary Wootters. Can we access a database both locally and privately? In Yael Kalai and Leonid Reyzin, editors, *Theory of Cryptography*, pages 662–693, Cham, 2017. Springer International Publishing.
- [Blo70] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426, 1970.
- [BT21] Alexandra Boldyreva and Tianxin Tang. Privacy-preserving approximate k-nearest-neighbors search that hides access, query and volume patterns. *PoPETS*, 2021.
- [CAD⁺23] Chloe Cachet, Sohaib Ahmad, Luke Demarest, Serena Riback, Ariel Hamlin, and Benjamin Fuller. Multi random projection inner product encryption, applications to proximity searchable encryption for the iris biometric. *Information and Computation*, 293:105059, 2023.
- [CCD⁺20] Hao Chen, Ilaria Chillotti, Yihe Dong, Oxana Poburinnaya, Ilya Razenshteyn, and M Sadegh Riazi. SANNs: Scaling up secure approximate k-Nearest neighbors search. In *USENIX Security*, pages 2111–2128, 2020.
- [CFR23] Anrin Chakraborti, Giulia Fanti, and Michael K Reiter. Distance-Aware private set intersection. In *USENIX Security*, pages 319–336, 2023.
- [CG97] Benny Chor and Niv Gilboa. Computationally private information retrieval (extended abstract). In *STOC, STOC '97*, page 304–313, New York, NY, USA, 1997. Association for Computing Machinery.
- [CGKO11] Reza Curtmola, Juan Garay, Seny Kamara, and Rafail Ostrovsky. Searchable symmetric encryption: Improved definitions and efficient constructions. *Journal of Computer Security*, 19:895–934, 01 2011.
- [CGPR15] David Cash, Paul Grubbs, Jason Perry, and Thomas Ristenpart. Leakage-abuse attacks against searchable encryption. In *CCS*, pages 668–679, 2015.
- [CHLR18] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. Labeled psi from fully homomorphic encryption with malicious security. In *CCS*, pages 1223–1237, 2018.
- [CHR17] Ran Canetti, Justin Holmgren, and Silas Richelson. Towards doubly efficient private information retrieval. In Yael Kalai and Leonid Reyzin, editors, *Theory of Cryptography*, pages 694–726, Cham, 2017. Springer International Publishing.
- [CKGS98] Benny Chor, Eyal Kushilevitz, Oded Goldreich, and Madhu Sudan. Private information retrieval. *J. ACM*, 45(6):965–981, nov 1998.
- [CLT16] Binyi Chen, Huijia Lin, and Stefano Tessaro. Oblivious parallel ram: improved efficiency and generic constructions. In *Theory of Cryptography Conference*, pages 205–234. Springer, 2016.
- [CT10] Emiliano De Cristofaro and Gene Tsudik. Practical private set intersection protocols with linear computational and bandwidth complexity. 2010. <https://eprint.iacr.org/2009/491>.

- [CWD⁺18] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018.
- [Dau05] John Daugman. Results from 200 billion iris cross-comparisons. 01 2005.
- [Dau09] John Daugman. How iris recognition works. In *The essential guide to image processing*, pages 715–739. Elsevier, 2009.
- [Dau14] John Daugman. 600 million citizens of India are now enrolled with biometric id.”. *SPIE newsroom*, 7, 2014.
- [DCW13] Changyu Dong, Liqun Chen, and Zikai Wen. When private set intersection meets big data: an efficient and scalable protocol. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 789–800, 2013.
- [DGXZ19] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, pages 4690–4699, 2019.
- [DSMRY09] Dana Dachman-Soled, Tal Malkin, Mariana Raykova, and Moti Yung. Efficient robust private set intersection. In *International Conference on Applied Cryptography and Network Security*, pages 125–142. Springer, 2009.
- [EFG⁺09] Zekeriya Erkin, Martin Franz, Jorge Guajardo, Stefan Katzenbeisser, Inald Lagendijk, and Tomas Toft. Privacy-preserving face recognition. In Ian Goldberg and Mikhail J. Atallah, editors, *Privacy Enhancing Technologies*, pages 235–253, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
- [FAKM14] Bin Fan, Dave G Andersen, Michael Kaminsky, and Michael D Mitzenmacher. Cuckoo filter: Practically better than bloom. In *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies*, pages 75–88, 2014.
- [FMC⁺15] Benjamin Fuller, Darby Mitchell, Robert Cunningham, Uri Blumenthal, Patrick Cable, Ariel Hamlin, Lauren Milechin, Mark Rabe, Nabil Schear, Richard Shay, et al. Security and privacy assurance research (spar) pilot final report. Technical report, MIT Lincoln Laboratory Lexington United States, 2015.
- [FMC⁺20] Francesca Falzon, Evangelia Anna Markatou, David Cash, Adam Rivkin, Jesse Stern, and Roberto Tamassia. Full database reconstruction in two dimensions. In *CCS*, pages 443–460, 2020.
- [FNP04] Michael J Freedman, Kobbi Nissim, and Benny Pinkas. Efficient private matching and set intersection. In *International conference on the theory and applications of cryptographic techniques*, pages 1–19. Springer, 2004.
- [Fou] Electronic Frontier Foundation. Mandatory national ids and biometric databases.
- [FP22] Francesca Falzon and Kenneth G. Paterson. An efficient query recovery attack against a graph encryption scheme. In Vijayalakshmi Atluri, Roberto Di Pietro, Christian D. Jensen, and Weizhi Meng, editors, *Computer Security – ESORICS 2022*, pages 325–345, Cham, 2022. Springer International Publishing.
- [FVK⁺15] Ben A Fisch, Binh Vo, Fernando Krell, Abishek Kumarasubramanian, Vladimir Kolesnikov, Tal Malkin, and Steven M Bellovin. Malicious-client security in blind seer: a scalable private dbms. In *2015 IEEE Symposium on Security and Privacy*, pages 395–410. IEEE, 2015.
- [FVY⁺17] Benjamin Fuller, Mayank Varia, Arkady Yerukhimovich, Emily Shen, Ariel Hamlin, Vijay Gadepally, Richard Shay, John Darby Mitchell, and Robert K Cunningham. SoK: Cryptographically protected database search. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 172–191. IEEE, 2017.
- [FWG⁺16] Zhangjie Fu, Xinle Wu, Chaowen Guan, Xingming Sun, and Kui Ren. Toward efficient multi-keyword fuzzy search over encrypted outsourced data with accuracy improvement. *IEEE Transactions on Information Forensics and Security*, 11(12):2706–2716, 2016.

- [Gen09] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *STOC*, pages 169–178, 2009.
- [GL20] Thomas Mueller Graf and Daniel Lemire. Xor filters: Faster and smaller than bloom and cuckoo filters. *Journal of Experimental Algorithmics (JEA)*, 25:1–16, 2020.
- [GLMP18] Paul Grubbs, Marie-Sarah Lacharité, Brice Minaud, and Kenneth G Paterson. Pump up the volume: Practical database reconstruction from volume leakage on range queries. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 315–331, 2018.
- [GMP16] Sanjam Garg, Payman Mohassel, and Charalampos Papamanthou. Tworam: efficient oblivious ram in two rounds with applications to searchable encryption. In *Annual International Cryptology Conference*, pages 563–592. Springer, 2016.
- [GO96] Oded Goldreich and Rafail Ostrovsky. Software protection and simulation on oblivious rams. *Journal of the ACM (JACM)*, 43(3):431–473, 1996.
- [Gol87] Oded Goldreich. Towards a theory of software protection and simulation by oblivious rams. In *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, pages 182–194, 1987.
- [GPAM⁺20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [GPP23] Zichen Gui, Kenneth G. Paterson, and Sikhar Patranabis. Rethinking searchable symmetric encryption. In *IEEE Security and Privacy*, 2023. <https://eprint.iacr.org/2021/879>.
- [GRGB⁺12] Javier Galbally, Arun Ross, Marta Gomez-Barrero, Julian Fierrez, and Javier Ortega-Garcia. From the iriscode to the iris: A new vulnerability of iris recognition systems. *Black Hat Briefings USA*, 1, 2012.
- [GRS22] Gayathri Garimella, Mike Rosulek, and Jaspal Singh. Structure-aware private set intersection, with applications to fuzzy matching. In *Annual International Cryptology Conference*, pages 323–352. Springer, 2022.
- [GSB⁺17] Paul Grubbs, Kevin Sekniqi, Vincent Bindschaedler, Muhammad Naveed, and Thomas Ristenpart. Leakage-abuse attacks against order-revealing encryption. In *Security and Privacy (SP), 2017 IEEE Symposium on*, pages 655–672. IEEE, 2017.
- [Hac18] Gabriel Hackebeil. Path oram python module, 2018.
- [HWKL18] Yi Huang, Adams Kong Wai-Kin, and Kwok-Yan Lam. From the perspective of CNN to adversarial iris images. In *BTAS*, pages 1–10. IEEE, 2018.
- [IKK12] Mohammad Saiful Islam, Mehmet Kuzu, and Murat Kantarcioglu. Access pattern disclosure on searchable encryption: ramification, attack and mitigation. In *NDSS*, volume 20, page 12. Citeseer, 2012.
- [IM98] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *STOC*, pages 604–613, 1998.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, page 448–456. JMLR.org, 2015.
- [JM19] Alexia Jolicoeur-Martineau. The relativistic discriminator: a key element missing from standard GAN. In *ICLR*, 2019.
- [KB15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [KE19] Evgenios M Kornaropoulos and Petros Efstathopoulos. The case of adversarial inputs for secure similarity approximation protocols. In *EuroS&P*, pages 247–262. IEEE, 2019.

- [KIK12] Mehmet Kuzu, Mohammad Saiful Islam, and Murat Kantarcioglu. Efficient similarity search over encrypted data. In *2012 IEEE 28th International Conference on Data Engineering*, pages 1156–1167. IEEE, 2012.
- [KKM⁺22] Seny Kamara, Abdelkarim Kati, Tarik Moataz, Thomas Schneider, Amos Treiber, and Michael Yonli. Sok: Cryptanalysis of encrypted search with leaker—a framework for leakage attack evaluation on real-world data. In *EuroS&P*, pages 90–108. IEEE, 2022.
- [KKNO16] Georgios Kellaris, George Kollios, Kobbi Nissim, and Adam O’Neill. Generic attacks on secure outsourced databases. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1329–1340, 2016.
- [KPT19] Evgenios M Kornaropoulos, Charalampos Papamanthou, and Roberto Tamassia. Data recovery on encrypted databases with k-nearest neighbor query leakage. In *IEEE S&P*, pages 1033–1050. IEEE, 2019.
- [KPT20a] Evgenios M Kornaropoulos, Charalampos Papamanthou, and Roberto Tamassia. The state of the uniform: attacks on encrypted databases beyond the uniform query distribution. In *IEEE S&P*, pages 1223–1240. IEEE, 2020.
- [KPT20b] Evgenios M. Kornaropoulos, Charalampos Papamanthou, and Roberto Tamassia. The state of the uniform: Attacks on encrypted databases beyond the uniform query distribution. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1223–1240. IEEE, 2020.
- [KYV⁺17] Naman Kohli, Daksha Yadav, Mayank Vatsa, Richa Singh, and Afzel Noore. Synthetic iris presentation attack using idcgan. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 674–680. IEEE, 2017.
- [LMW23] Wei-Kai Lin, Ethan Mook, and Daniel Wichs. Doubly efficient private information retrieval and fully homomorphic ram computation from ring lwe. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing, STOC 2023*, page 595–608, New York, NY, USA, 2023. Association for Computing Machinery.
- [LMWY20] Kasper Green Larsen, Tal Malkin, Omri Weinstein, and Kevin Yeo. Lower bounds for oblivious near-neighbor search. In *SODA*, pages 1116–1134. SIAM, 2020.
- [LN⁺96] Michael Lamoureux, Bradford Nickerson, et al. On the equivalence of b-trees and deterministic skip list. 1996.
- [LPW⁺20] Qin Liu, Yu Peng, Jie Wu, Tian Wang, and Guojun Wang. Secure multi-keyword fuzzy searches with enhanced service quality in cloud computing. *IEEE Transactions on Network and Service Management*, 18(2):2046–2062, 2020.
- [LWW⁺10] Jin Li, Qian Wang, Cong Wang, Ning Cao, Kui Ren, and Wenjing Lou. Fuzzy keyword search over encrypted data in cloud computing. In *INFOCOM, 2010 Proceedings IEEE*, pages 1–5. IEEE, 2010.
- [MCYJ18] Guangcan Mai, Kai Cao, Pong C Yuen, and Anil K Jain. On the reconstruction of face images from deep face templates. *IEEE transactions on pattern analysis and machine intelligence*, 41(5):1188–1202, 2018.
- [MHN13] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3, 2013.
- [MT19] Evangelia Anna Markatou and Roberto Tamassia. Full database reconstruction with access and search pattern leakage. In *International Conference on Information Security*, pages 25–43. Springer, 2019.
- [OPJM10] Margarita Osadchy, Benny Pinkas, Ayman Jarrous, and Boaz Moskovich. SCiFI - a system for secure face identification. In *2010 IEEE Symposium on Security and Privacy*, pages 239–254, 2010.

- [OPPW24] Hiroki Okada, Rachel Player, Simon Pohmann, and Christian Weinert. Towards practical doubly-efficient private information retrieval. In *Financial Cryptography*, 2024.
- [PBF⁺08] P Jonathon Phillips, Kevin W Bowyer, Patrick J Flynn, Xiaomei Liu, and W Todd Scruggs. The iris challenge evaluation 2005. In *BTAS*, pages 1–8. IEEE, 2008.
- [PKV⁺14] Vasilis Pappas, Fernando Krell, Binh Vo, Vladimir Kolesnikov, Tal Malkin, Seung Geol Choi, Wesley George, Angelos Keromytis, and Steve Bellovin. Blind seer: A scalable private dbms. In *IEEE S&P*, pages 359–374. IEEE, 2014.
- [pri24] Implementation of private eyes. <https://github.com/PrivateEyes-paper/PrivateEyes>, 2024.
- [PSO⁺09] P Jonathon Phillips, W Todd Scruggs, Alice J O’Toole, Patrick J Flynn, Kevin W Bowyer, Cathy L Schott, and Matthew Sharpe. FRVT 2006 and ICE 2006 large-scale experimental results. *IEEE transactions on pattern analysis and machine intelligence*, 32(5):831–846, 2009.
- [Pug90] William Pugh. Skip lists: a probabilistic alternative to balanced trees. *Communications of the ACM*, 33(6):668–676, 1990.
- [RS21] Peter Rindal and Phillipp Schoppmann. Vole-psi: Fast oprf and circuit-psi from vector-ole. In Anne Canteaut and François-Xavier Standaert, editors, *Advances in Cryptology – EUROCRYPT 2021*, pages 901–930, Cham, 2021. Springer International Publishing.
- [SDDN19] Sobhan Soleymani, Ali Dabouei, Jeremy Dawson, and Nasser M Nasrabadi. Adversarial examples to fool iris recognition systems. In *2019 International Conference on Biometrics (ICB)*, pages 1–8. IEEE, 2019.
- [SDS⁺18] Emil Stefanov, Marten Van Dijk, Elaine Shi, T-H Hubert Chan, Christopher Fletcher, Ling Ren, Xiangyao Yu, and Srinivas Devadas. Path oram: an extremely simple oblivious ram protocol. *Journal of the ACM (JACM)*, 65(4):1–26, 2018.
- [SSF19] Sailesh Simhadri, James Steel, and Benjamin Fuller. Cryptographic authentication from the iris. In *International Conference on Information Security*, pages 465–485. Springer, 2019.
- [SWP00] Dawn Xiaoding Song, David Wagner, and Adrian Perrig. Practical techniques for searches on encrypted data. In *IEEE S&P*, pages 44–55. IEEE, 2000.
- [UCK⁺21] Erkam Uzun, Simon P Chung, Vladimir Kolesnikov, Alexandra Boldyreva, and Wenke Lee. Fuzzy labeled private set intersection with applications to private real-time biometric search. In *USENIX Security*, pages 911–928, 2021.
- [VS11] Shreyas Venugopalan and Marios Savvides. How to generate spoofed irises from an iris code template. *IEEE Transactions on Information Forensics and Security*, 6(2):385–395, 2011.
- [WLD⁺17] Guofeng Wang, Chuanyi Liu, Yingfei Dong, Hezhong Pan, Peiyi Han, and Binxing Fang. Query recovery attacks on searchable encryption based on partial knowledge. In *International Conference on Security and Privacy in Communication Systems*, pages 530–549. Springer, 2017.
- [WMT⁺13] Jianfeng Wang, Hua Ma, Qiang Tang, Jin Li, Hui Zhu, Siqi Ma, and Xiaofeng Chen. Efficient verifiable fuzzy keyword search over encrypted data in cloud computing. *Comput. Sci. Inf. Syst.*, 10(2):667–684, 2013.
- [WNL⁺14] Xiao Shaun Wang, Kartik Nayak, Chang Liu, TH Hubert Chan, Elaine Shi, Emil Stefanov, and Yan Huang. Oblivious data structures. In *CCS*, pages 215–226, 2014.
- [WST12] Peter Williams, Radu Sion, and Alin Tomescu. Privatefs: A parallel oblivious file system. In *CCS*, pages 977–988, 2012.
- [WYLH14] Bing Wang, Shucheng Yu, Wenjing Lou, and Y Thomas Hou. Privacy-preserving multi-keyword fuzzy search over encrypted data in the cloud. In *IEEE INFOCOM 2014-IEEE conference on computer communications*, pages 2112–2120. IEEE, 2014.

Layer	OutputSize	Kernels
Input	1x128	-
G1	64	64
G2	128	128
G3	64	64
G4	128	128
G5	1024	1024
D1	16	16
D2	128	128
D3	1	1

Table 4: Generator architecture

[YCR19] Shivangi Yadav, Cunjian Chen, and Arun Ross. Synthesizing iris images using RaSGAN with application in presentation attack detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.

A Synthetic Data Generation

We now describe the neural network used to produce our synthetic templates. Our synthetic templates are built using a generative adversarial network or GAN. A GAN trains two networks in competition, a generator which should produce synthetic irises and a discriminator which classifies irises as real or synthetic. Yadav et al. [YCR19] uses RaSGAN (relativistic average standard GAN) [JM19] to generate synthetic irises for the purpose of studying their effects on presentation attack detection (PAD) algorithms. Irises from the RaSGAN perform well against PAD and follow real iris statistics well. Kohli et al. [KYV⁺17] use the DCGAN architecture to generate synthetic irises. Synthetic irises can be viewed as irises that must closely resemble bonafide irises as discussed in [YCR19, KYV⁺17].

We follow the approach of RESIST [AF20] which takes inspiration from synthetic data generation to invert iris templates into realistic looking images. We denote the network as SYNTH. Training uses the ND-0405 dataset [BF16]. In a GAN formulation, noise is sampled from a multivariate normal distribution (\mathbb{P}_z) with a mean of 0 and a variance of 1. The generator converts this noise vector into a synthetic template. Let \mathbb{P}_y denote the distribution of real templates and $\mathbb{P}_{\hat{y}}$ denote the distribution of synthetic templates. We use a recently proposed relativistic average discriminator [JM19] as our discriminator. To build up to the relativistic discriminator we first start with the original GAN loss functions:

$$L(D) = -\mathbb{E}_{y \sim \mathbb{P}_y} [\log(D(y))] - \mathbb{E}_{\hat{y} \sim \mathbb{P}_{\hat{y}}} [\log(1 - D(\hat{y}))]$$

$$L(G) = -\mathbb{E}_{\hat{y} \sim \mathbb{P}_{\hat{y}}} [\log(D(\hat{y}))].$$

$L(D)$ is called the discriminator loss and $L(G)$ is called the generator loss, y is an actual template and \hat{y} is a synthetic template generated by the generator.

The $L(D)$ and $L(G)$ losses are minimized using gradient descent. The generator and discriminator play a zero sum game. The generator weights are updated based on how good its synthetic irises are while the discriminator weights are updated on how well it differentiates between real and synthetic templates.

The last layer of the generator is the hyperbolic tangent (tanh) with output ranging from -1 to 1, this is done to generate binary synthetic templates by substituting 0,1 with -1,1. The real templates are also converted to -1,1 for conformity.

Architecture SYNTH architecture is a small neural network having only dense (fully connected) layers as shown in Table 4 where G_x are generator layers and D_x are discriminator layers. Each layer is followed by a LeakyReLU [MHN13] activation and a batch normalization [IS15] layer. The last layers of both sub-networks are unique, the generator has a tanh activation while the discriminator has a Sigmoid activation.

Training SYNTH is trained in two stages. First, the generator produces a synthetic template and second, the discriminator outputs how real this synthetic template is. This training is done till convergence of the weights of both networks. Both training stages use the Adam optimizer [KB15]. We randomly flip 2% bits of a real template as noise to aid in network convergence. The networks is trained over 100 epochs. Each epoch having 100 steps. Each step updates weights once by either 1) discriminating a pair of vectors or 2) generating a single synthetic template.

Once trained the SYNTH network can produce an unbounded number of distinct templates. We described how to produce different readings from the same template in Section 5.