

HEBI: Homomorphically Encrypted Biometric Indexing

Pia Bauspieß^{*†}, Marcel Grimmer^{*}, Cecilie Fougner^{*}, Damien Le Vasseur[‡], Thomas Thaulow Stöcklin^{*}, Christian Rathgeb[†], Jascha Kolberg[†], Anamaria Costache^{*}, Christoph Busch^{*†}

^{*}NTNU – Norwegian University of Science and Technology

[†]Hochschule Darmstadt – da/sec Biometrics and Security Research Group

[‡]ENSIIE – National School of Computer Science for Industry and Business

{pia.bauspiess, marceg, anamaria.costache, christoph.busch}@ntnu.no

Abstract

Biometric data stored in automated recognition systems are at risk of attacks. This is particularly true for large-scale biometric identification systems, where the reference database is often accessed remotely. A popular approach for the protection of the stored templates is homomorphic encryption, which grants privacy protection while maintaining the biometric performance of the unprotected system. However, it introduces a significant computational overhead that can render identification transactions infeasible. To reduce this workload, biometric indexing in the encrypted domain has become a recent research interest. In this work, we show that in such schemes, auxiliary indexing data can leak additional privacy-sensitive information that violate standardized requirements for biometric template protection. In response to this leakage, we propose a novel framework HEBI that protects biometric indexing approaches at a post-quantum security level while requiring a computational effort of only 0.12 milliseconds per cluster.

1. Introduction

Biometric data allow for an irrevocable identification of individuals over several decades [28]. Therefore, biometric data need to be considered sensitive data requiring long-term protection, even more so than passwords or authentication tokens that can be exchanged upon a security breach. To ensure this protection, the ISO/IEC 24745 standard on biometric information protection [25] defines the following requirements: *i) unlinkability*, two protected templates stored in different applications cannot be linked to the same subject, *ii) renewability*, new templates can be created from the same source if the previously stored reference was leaked without the need to re-enrol, and *iii) irreversibility*, it is impossible to reconstruct original samples given only protected templates. Furthermore, both the computational and biometric performance (i.e., accuracy) of the unprotected

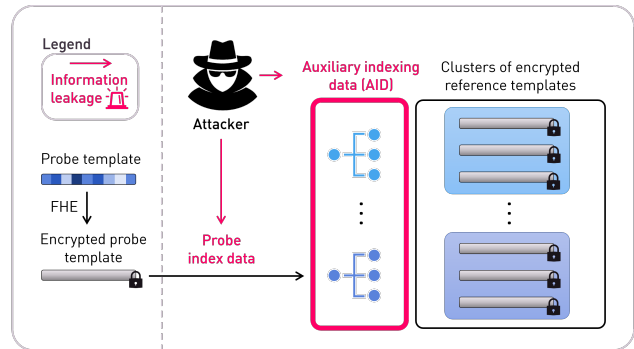


Figure 1. Biometric information leakage in indexing schemes on encrypted reference databases: an attacker can observe privacy-sensitive information from indexing data, e.g., soft-biometric attributes such as gender of the probe and reference subjects.

system should be preserved.

In biometric identification, where a $1 : N$ search against a large database is performed, biometric templates are at particular risk as reference databases are maintained for long time spans. For example, this is true for criminal databases held by law enforcement agencies or for national citizen registration [42]. In addition, these databases are static targets of attack, as their large storage requirements do not allow for agile changes to their physical security.

Recently, biometric identification protected through *Fully Homomorphic Encryption* (FHE) has been explored to mitigate these security risks [4, 14, 19]. While this approach grants cryptographically sound protection of the biometric templates, it comes with a significant overhead in computational workload. For large-scale databases, workload reduction strategies need to be applied to achieve practical biometric identification systems. Workload reduction strategies have been categorized into two main classes [15]: feature transformation and preselection. Preselection approaches offer a significant speed-up through selecting a smaller subset of the enrolment database that contains the reference identifier with high probability. Using an index

string i common to a subset of enrolled references \mathcal{C}_i , preselection can be achieved in $\mathcal{O}(1)$ and is therefore efficient.

However, a key challenge with these approaches is the continuous protection of data subject privacy under preselection, i.e., ensuring that the preselection procedure and its outcome do not reveal any information about the underlying subject, or infringe on the unlinkability of the system. This vulnerability is depicted in Figure 1. It is important to note that the encryption of the feature vectors alone is not sufficient to fulfil this requirement, as the preselection algorithm can reveal additional information about the enrolled subjects, e.g., soft-biometric characteristics such as the gender or ethnicity of the probe and reference subjects.

The risk of information leakage shown in Figure 1 is particularly high when biometric indexing is based on similarity measures between the enrolled subjects, e.g., in feature-based clustering approaches. These similarity measures contained in the *Auxiliary Indexing Data (AID)* can potentially reveal sensitive information about the preselected subset such as their shared soft-biometric characteristics. For sound privacy protection in the sense of ISO/IEC 24745 [25], this information needs to be obscured in addition to the protection of the feature vectors.

To mitigate the privacy leakage in biometric indexing, we therefore propose a novel protocol HEBI that can be applied to indexing approaches in the encrypted domain. The key contributions of our work are as follows:

- **Privacy analysis.** To illustrate the significant risks that come with the use of unprotected AID, we give a privacy analysis of existing approaches. We show that we were able to reconstruct the gender and ethnicity of enrolled subjects based only on AID, which must be considered a severe security risk.
- **Formalization of information leakage.** Further, we give a formalization of information leakage in biometric indexing that indicates that such a leakage exists in arbitrary biometric indexing schemes. We use this formalization as further motivation for our work, in addition to the experimental analysis.
- **The novel HEBI protocol.** As our main contribution, we present the HEBI protocol for secure biometric indexing in the encrypted domain. Through the use of lattice-based cryptography [5, 9], our protocol provides post-quantum security in storage, preselection and comparison. We give an experimental evaluation that shows that HEBI can be applied in real-world operational systems at a cost of only 0.12 additional milliseconds for the the post-quantum secure retrieval compared to unprotected preselection systems. At the same time, the biometric performance of the underlying indexing approach is not impacted by the applied cryptographic protection mechanisms.

- **Security analysis.** We provide a comprehensive security analysis of our protocol and show how it mitigates the flaws of unprotected approaches, thus giving full post-quantum security to biometric data under preselection.

The remainder of this article is structured as follows: Section 2 discusses works that are closely related to ours, before we analyse of the privacy leakage in a previously proposed privacy-preserving biometric indexing scheme. In Section 4, we introduce more technical cryptographic background information. From this, we present our novel HEBI protocol in Section 5 that alleviates the presented privacy risks. Section 6 gives experimental results and a security analysis. Finally, we draw our conclusions in Section 7.

2. Related Work

Workload reduction in homomorphically encrypted biometric identification systems has recently been achieved with post-quantum security [4, 19]. However, both of these works were only based on feature transformation, such that an exhaustive search requiring a linearly increasing costs remains. It is important to note that our HEBI protocol can integrate such feature transformation approaches seamlessly and therefore allows for further improvements in large-scale biometric identification systems.

The cryptographic concept of homomorphic search has previously been applied to biometric identification in [46]. In their work, the authors use the search scheme as a replacement for FHE rather than an additional protection layer for the preselection step. In order to realize homomorphic search on the feature vectors directly, strong statistical assumptions about the feature representation are required, which do not generalize over different modalities. Another recent work [45] applied homomorphic search for biometric authentication instead of identification. Most recently, [3] applied homomorphic search for preselection on an encrypted reference database. However, our HEBI protocol differs non-trivially from the proposal in [3] in several aspects. Firstly, the work by [3] can only be considered as proof-of-concept, as a handcrafted preselection approach is utilized in their work, which underlies the unrealistic assumption of perfect ground truth. In comparison, HEBI is designed for real-world indexing approaches that allow for a meaningful analysis of the overall biometric performance. Secondly, [3] apply a binning approach that does not trivially generalize to other application scenarios apart from their own, while HEBI enables efficient and secure cluster generation independent of the indexing algorithm. Finally, our work offers an extensive analysis of the risk of preselection independent of the concrete indexing approach and shows how to mitigate these risks in a universal approach.

The application of unprotected biometric indexing to

biometric identification [13, 24, 33, 34, 35, 38, 41, 43] will be discussed at length in the following Section. These are the schemes our HEBI protocol improves upon through an additional layer of protection during the preselection step. Notably, the choice of protection mechanism for the reference database is independent of the HEBI preselection protocol, though we adhere to FHE-based protection in our work. In addition, HEBI does not impair the originally given biometric performance of the above works.

3. Privacy Analysis of Biometric Indexing

Biometric indexing as depicted in Figure 1 has been applied in a number of recent research works, among others [13, 24, 33, 34, 35, 38, 41, 43]. In this Section, we give further intuition to the privacy implications of such approaches through probability theory.

3.1. Formal Model

In this analysis, we investigate the relation between the enrolled reference feature vectors $\{r_j\}_{j=0}^{N-1}$ for a number of reference subjects N and the auxiliary indexing data (AID) represented by index strings $\{i_k\}_{k=0}^K$, where K denotes the number of clusters or index strings in the given scheme. We define that every reference feature vector r_j is assigned one and only one index string i_k , while one index string i_k clusters several references (i.e., $K < N$). Upon an identification transaction, a probe feature vector p is extracted from a presented probe sample, and the corresponding index i_k is determined. Then, only the reference features vectors associated with i_k are compared to p in the encrypted domain.

For the formalization of privacy leakage in such indexing schemes, we utilize the information-theoretic concept of *mutual information* $I(X; Y)$, which is defined as

$$I(X; Y) = D_{KL}(P_{(X,Y)} || P_X \otimes P_Y), \quad (1)$$

where X and Y are random variables and D_{KL} denotes the Kullback–Leibler divergence [11]. The mutual information can further be expressed in terms of entropy [40]:

$$I(X; Y) = H(X) - H(X|Y), \quad (2)$$

where $H(X)$ is the marginal entropy of X and $H(X|Y)$ is the conditional entropy of X given Y . Let $\{X\}_j$ be the variable family that represents the reference feature vectors and $\{Y\}_k$ be the variable family that represents the index strings. In a meaningful indexing scheme, it holds that

$$I(X_k, Y_{i_k}) > I(X_k, Y_{i_m}), \quad (3)$$

i.e., the mutual information between the reference feature vector r_j associated with index string i_k should be greater than the mutual information between the same reference feature vector r_j and a different cluster associated with an

index string i_m . Otherwise, r_j would be associated with i_m instead. From Equation 3, it follows that $H(X_k|Y_{i_m}) > H(X_k|Y_{i_k})$. As $H(X_k|Y_{i_m})$ cannot be smaller than 0, it follows that $H(X_k|Y_{i_k}) > 0$. At the same time, the similarity of index strings does not correspond to the full feature vectors, which would yield no advantage over an exhaustive identification search. Therefore,

$$H(X_k) > H(X_k|Y_{i_k}) > 0, \quad (4)$$

and consequently,

$$I(X_k; Y_{i_k}) = H(X_k) - H(X_k|Y_{i_k}) > 0, \quad (5)$$

meaning that there is mutual information contained between the feature vectors and index strings. This mutual information defines the leakage of biometric information, which allows for attacks on the probe and reference subjects that can violate their privacy. Indeed, it has been shown that auxiliary data in biometric systems can lead to privacy risks in other applications, e.g., biometric cryptosystems [39]. However, we emphasize that our formal model is not intended to be used as a concrete metric, as mutual information is hard to calculate precisely. Instead, it serves as a logical argument for the existence of privacy leakage in AID.

More empirically, index strings are commonly constructed such that they allow for a clustering of the reference feature vectors based on a more general measure of similarity than the exact comparison between feature vectors. In some approaches [13, 34, 38], the index strings are even derived from the feature vectors directly, representing a down-sampled representation of one or more feature vectors. In the following, we show how to concretely extract privacy-sensitive information from such representations.

3.2. Case Study

To illustrate the risks of soft-biometric leakage in biometric indexing in a case study, we analyze the recent work of [34], which is one of the works relying on unprotected index strings discussed above.

In their work, the authors generate a look-up table of short binary strings, or *stable hashes*, which represent distinct clusters of reference templates. They present different methods of obtaining these stable hashes, all of which are based on the feature representations of the enrolled references. In our evaluation, we focus in the first of their proposed approaches, which is the established k-means clustering technique [32]. During the enrolment phase, the clustering algorithm is trained on the enrolment database, which is subsequently encrypted using FHE. The protected templates are stored in the database alongside the look-up table of stable hashes, which in the case of k-means clustering are a binary representation of the cluster centers, or centroids. Upon an identification transaction, the distance of the probe

feature vector to all centroids is calculated, and the closest centroid is determined to be the probe stable hash. Then, the reference subjects with the same stable hash are extracted from the enrolment database, a homomorphic comparison of the encrypted probe feature vector against the encrypted references is computed, and the decision is revealed to the client that initiated the transaction [34].

The advantage of this indexing approach is the error-correcting capability of the clustering approach, which allows for an exact comparison of the stable hashes and is therefore very efficient. The retrieval cost of the look-up operation is constant at $\mathcal{O}(1)$ and can be considered negligible compared to the cost of the homomorphic operations. Furthermore, the low preselection error even on challenging data sets makes the approach in [34] attractive.

However, the vulnerability of the approach with regard to the reference subjects’ privacy lies in the stable hash look-up table, which is stored alongside the enrolment database. As argued above, it can be expected that the stable hashes encode information about the probe and reference subjects to some degree, which could be privacy-sensitive information. For example, soft-biometric similarities to the subjects in one cluster could be revealed, which would constitute a violation of ISO/IEC 24745 [25]. Disclosure of soft-biometric data related to the ethnic origin is a breach of the European Union’s General Data Protection Regulation [20].

To confirm our hypothesis, we conducted an experimental evaluation of the privacy leakage in the system presented in [34]. For this evaluation, we selected 3,165 samples of 533 subjects from the Face Recognition Grand Challenge (FRGC) database [36] that are compliant with the International Civil Aviation Organization’s face image quality requirements for machine-readable travel documents. The code for the stable hash generation from k-means clustering [32] has been provided by the authors to facilitate the reproducibility of their results. In terms of parameters, we followed the original work with $P = 1$ subspaces and $K = 64$.

Figures 2 and 3 show the distribution of ethnicity and gender of the 64 clusters. For this analysis, ground truth labels for the image samples were hand-annotated, such that a high accuracy in the labelling can be assumed compared to soft-biometric feature extractors [17]. From the visualization of the distributions, it becomes evident that there exists pooling of soft-biometric characteristics within both dimensions of ethnicity and gender. This can be for example observed in clusters clusters 36, 38 and 39, which exclusively contain female subjects, while clusters 3, 9 and 11 only contain male subjects. Similarly, clusters 10 and 11 exclusively contain Asian subjects, while clusters 21, 22, 23, 42, 46 and 63 only contain Caucasian subjects. While these characteristics are not perfectly separated over all clusters, it is particularly concerning that the clustering effectively exposes underrepresented subgroups. A prominent example is

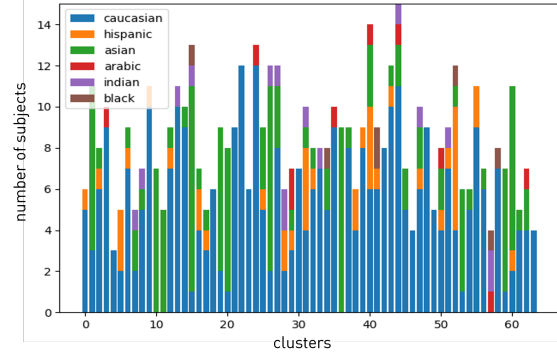


Figure 2. Distribution of ethnicities over the clusters derived from an ICAO-compliant subset of the FRGC database [36].

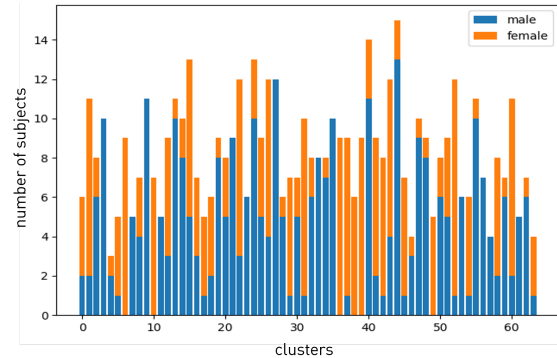
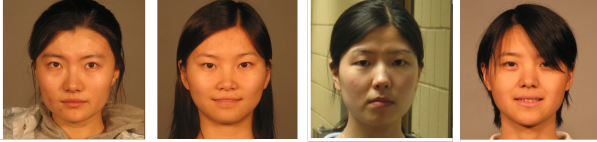


Figure 3. Distribution of genders over the clusters derived from an ICAO-compliant subset of the FRGC database [36].

cluster 10, which contains only female Asian subjects. An attacker observing the stable hash corresponding to cluster 10 can therefore with high probability deduce the gender and ethnicity of the probe subject and the reference subjects stored alongside that stable hash.

To extend our analysis, we further evaluated a synthetic face image generation from the centroids to approximate the average features of the subjects in the clusters and their similarity to the synthetic approximation for the respective cluster. We leveraged the StyleGAN3 generator [26] pre-trained on the FFHQ database [27] that includes more than 70,000 face images with diverse ethnicities, gender labels, and other facial characteristics. To reconstruct latent representations and subsequently derived representative face images from each stable hash (s), we trained a fully connected neural network (\mathcal{M}) that maps each stable hash into the semantic manifold of the StyleGAN3 intermediate latent space. We froze the generator (\mathcal{G}) weights during training to preserve its capability to generate photo-realistic face images. Further, we applied a simple *mean squared error* loss function to minimize the difference between the reconstructed face images $\hat{x} = \mathcal{G}(\mathcal{M}(s))$ to the randomly drawn face images x of their corresponding cluster.

For this experiment, the FRGCv2 training subset has been reduced such that each stable hash is assigned with only one face image per identity. This setting prevents



(a) Bona fide samples of subjects from cluster 10.



(b) Left to right: Reconstructions based on untrained, trained on cluster 10 only, trained on all clusters, trained on 70% of all clusters (excluding cluster 10) StyleGAN approximations.

Figure 4. Comparison of bona fide FRGCv2 samples of cluster 10 and StlyeGAN presentation attack approximations of cluster 10.

the mapping network from oscillating due to the high intra-subject variance. For the optimization of \mathcal{M} , the StyleGAN3 truncation factor was set to 0.75, enabling the generation of face images with stable quality. We adopted the Adam optimizer settings from [27] and increased the learning rate to 0.01 to accelerate the training process. The results of this evaluation are shown in Figure 4.

In Figure 4, the reconstructed latent representations corresponding to cluster 10 are depicted alongside a selection of bona fide sample images from that cluster, which contains only female Asian subjects. The reconstructed images are based on incrementally scarce training data to show that our GAN-based approach generalizes even in an open-set scenario. The closest approximation has been trained on cluster 10 alone, and cannot be considered a realistic attack. Both the closed-set and the open-set training scenario excluding cluster 10 continue however to show significant similarities to the original identities. Most importantly, the soft-biometric characteristics of gender and ethnicity are preserved. A breach of the latter in particular constitutes a GDPR [20] violation and must be prevented.

To conclude this analysis, significant privacy leakage has been found in the indexing approach by [34]. However, the overall indexing scheme is of high relevance to the problem of workload reduction for large-scale biometric identification, as it benefits from a high biometric performance and is therefore desirable to apply.

Looking towards the cryptographic protection of indexing approaches such as [13, 34, 38], the component of the index string that allows for the privacy leakage is their deterministic nature, i.e., in the case of [34], similar feature vectors will always be mapped to the same stable hash. In the remainder of this paper, we are therefore proposing a transformation of this deterministic preselection approach to a non-deterministic preselection, where similar feature vectors are mapped to randomized outputs that look

indistinguishable to an attacker. At the same time, they allow for the correct retrieval of the corresponding reference subjects, such that the biometric performance of the indexing approach is not impacted.

4. Preliminaries

4.1. Fully Homomorphic Encryption

Homomorphic encryption describes a cryptographic technique that allows for the evaluation of functions on encrypted data. More precisely, we call a public-key encryption scheme homomorphic if

$$\text{Enc}(pk_H, x \odot y) = \text{Enc}(pk_H, x) \odot \text{Enc}(pk_H, y). \quad (6)$$

More recently, *Fully Homomorphic Encryption* (FHE) has become practical for application in certain use cases. Following the groundbreaking work by Gentry [22], different schemes have established themselves with respect to their different properties. One of these is the CKKS [9] scheme, which provides the useful advantage of computing on high-precision approximations of floating point numbers directly, where other schemes require integer quantisation [8, 21] or binarisation [10]. In terms of the encrypted comparison of biometric feature vectors, this means that the underlying data does not need to be altered, and no information from the biometric comparison is lost. Therefore, the computations on encrypted templates correspond directly to computations on the unprotected templates, and the biometric performance remains unimpaired.

The security of many FHE schemes, including CKKS, is based on the Ring-Learning with Errors (R-LWE) problem, which is assumed to be secure against attacks implemented on a quantum computer [31]. These cryptosystems therefore provide a high level of protection to the biometric data, and in particular, long-term protection over several decades according to the current basis of knowledge and expectations in the field of cryptography [1].

4.2. Public-Key Encryption with Keyword Search

In addition to the protection of the feature vectors, the privacy analysis in Section 3 has shown that the indexing and retrieval during the preselection process requires additional protection. A recent work on face identification [3] has proposed the use of *Public-Key Encryption with Keyword Search* (PEKS) for the protection of semantic soft-biometric keywords. In this work, we apply this technique to generic biometric indexing approaches.

The cryptographic basis of PEKS lies in *Identity-Based Encryption* (IBE), which was first introduced by Boneh and Franklin in 2001 [7]. Building on this idea, PEKS was proposed as a means of creating ciphertexts for specific semantic keywords instead of identities [6]. In the typical application scenario, a PEKS scheme is used to create an en-

ryption of a keyword together with a corresponding trapdoor. This pair of cryptographic objects can be subjected to a publicly available test function which reveals no information except for the binary decision outcome of the similarity of the underlying keyword of the ciphertext and trapdoor.

A PEKS scheme [5] is defined as a tuple of four algorithms $\text{PEKS} = (\text{KeyGen}, \text{PEKS}, \text{Trapdoor}, \text{Test})$:

- $(pk_S, sk_S) \leftarrow \text{KeyGen}(1^k)$: On the input of the security parameter k , this algorithm outputs the public and secret key pair (pk_S, sk_S) .
- $s_w \leftarrow \text{PEKS}(pk_S, w)$: On the input of the public key pk_S and a keyword $w \in \{0, 1\}^*$, this algorithm outputs a searchable ciphertext s_w .
- $t_w \leftarrow \text{Trapdoor}(sk_S, w)$: On the input of a secret key sk_S and a keyword $w \in \{0, 1\}^*$, this algorithm outputs a trapdoor t_w .
- $b \leftarrow \text{Test}(t_w, s_w)$: On the input of a trapdoor $t_w = \text{Trapdoor}(sk_S, w')$ and a searchable ciphertext $s_w = \text{PEKS}(pk_S, w)$, this algorithm outputs a bit $b = 1$ if $w = w'$, and $b = 0$ otherwise.

More recently, PEKS has been implemented based on lattice-based IBE [18] to create lattice-based PEKS [5]. Compared to the original construction, lattice-based PEKS has high computational efficiency and provides post-quantum security through R-LWE [31]. As an important property to the application in this work, PEKS ciphertexts are constructed using a random component, yielding non-deterministic encryption. In the following Section, we will detail how this property ensures privacy protection when applied to biometric indexing.

5. The HEBI Protocol

In this Section, we present our novel HEBI protocol for biometric indexing in the encrypted domain. The protocol can be applied to any existing biometric indexing approach that clusters enrolment biometric references to prevent the leakage of sensitive information about the data subjects.

5.1. Setting

The HEBI protocol is executed between three parties: A *client* device, a *Database Server* (DS) and a *Trusted Third Party* service (TTP). All three parties are considered in the semi-honest security model, where they may aim to gain information about the data they are exchanging, but are not assumed to deviate from the given protocol. This is an established security assumption in remote biometric authentication [23, 29, 44].

5.2. Enrolment

During the enrolment phase, two separate setup operations are performed: initialisation of the encrypted indexing algorithm and the homomorphic encryption of the enrolment database.

The indexing algorithm is assumed to require some pre-computation on an unencrypted enrolment database [34]. In our protocol, we allow for this precomputation to be conducted during an offline phase prior to the deployment of the system, where the unprotected templates are not exposed to potential attacks. As a result of the indexing algorithm, each biometric reference r will be assigned an index, or cluster, i which can be of arbitrary data representation. If the clustering algorithm does not produce balanced clusters, i.e., the number of subjects per cluster is not consistent, the clusters are padded with random feature vectors to be of equal size.

Once the clusters have been established, the PEKS framework can be initialised. First, TTP generates a number of random PEKS keywords $\{w_i \mid 0 \leq i \leq K - 1\}$, where K is the total number of clusters, and fixes a mapping M between the clusters and PEKS keywords, which is made publicly available. Note that the mapping M of clusters to PEKS keywords must be indicated by the clusters' (arbitrarily assigned) order instead of the semantic index string i that could potentially reveal privacy-sensitive information. By making the mapping publicly available, the PEKS keywords do not act as additional secret keys in the system.

From the PEKS keywords, two look-up tables are generated. At TTP, a trapdoor $t_p \leftarrow \text{Trapdoor}(sk_S, w_i)$ is computed and stored for every cluster using the PEKS secret key sk_S . At DS, a mapping of encrypted references to clusters is stored, again based on any order of the clusters without using the index i as the identifier. An overview of the look-up tables is given in Figure 6.

For the setup of the encrypted enrolment database, TTP generates and stores a key pair of the homomorphic encryption scheme (sk_H, pk_H) and makes pk_H available to the client and DS. For a reference feature vector r , the client can enrol a data subject by computing $c_r \leftarrow \text{Enc}(pk_H, r)$ and sending c_r encrypted biometric reference to DS. Since the assignment of subjects to clusters is initially fixed, coefficient packing can be applied to facilitate further workload reduction [4].

5.3. Identification

During an identification transaction in HEBI, the client captures a probe sample and obtains its feature representation p . The client determines the index i_p of the probe with respect to the applied indexing algorithm and uses the public mapping M to determine the corresponding PEKS keyword w_p . Using the public key pk_H of the HE scheme, the client encrypts the probe feature vector by computing $c_p \leftarrow \text{Enc}(pk_H, p)$. It further computes the encrypted probe index $s_p \leftarrow \text{PEKS}(pk_S, w_p)$, and sends c_p and s_p to DS, which forwards s_p to TTP.

Upon receiving s_p , TTP determines the corresponding trapdoor t_i for which $\text{Test}(t_i, s_p) = 1$ holds true. Using the look-up table mapping trapdoors to clusters (see Figure

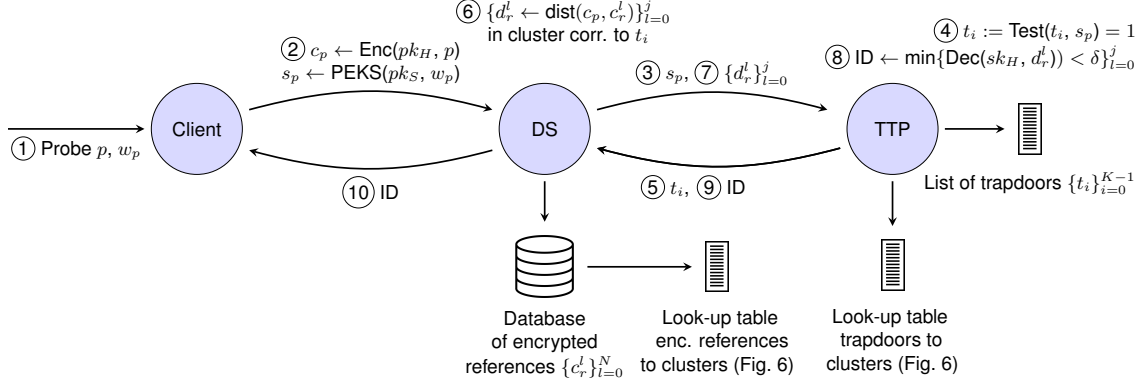


Figure 5. Identification transaction for encrypted preselection in the HEBI protocol.

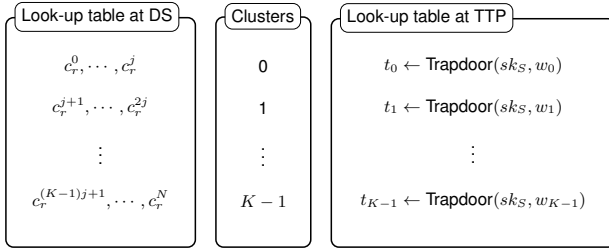


Figure 6. HEBI look-up tables generated at enrolment.

6), TTP sends the cluster identifier to DS, where the homomorphic comparisons are computed between the encrypted probe c_p and the encrypted references $\{c_r^l\}_{l=0}^j$ in the cluster corresponding to t_i . The encrypted comparison scores are sent to TTP, which decrypts them and determines the identification outcome, which is forwarded to the client. Note that throughout this transaction, DS and TTP do not have access to unprotected feature vectors or the index strings i that could reveal sensitive information. An overview of an identification transaction is given in Figure 5.

Our HEBI protocol can be seen as an independent layer of protection to arbitrary indexing schemes. Furthermore, it can also be combined with interchangeable template protection approaches for the feature vectors themselves, e.g., different FHE schemes or irreversible feature transformations. It is therefore versatile in its application and can be considered for applications beyond face recognition.

6. Experimental Evaluation

To show the practicality of our HEBI protocol, we give an experimental evaluation for the application to stable hashes [34]. By applying the additional layer of security, the privacy concerns outlined in Section 3 will be mitigated.

The experiments were conducted on the same subset of 533 subjects the FRGCv2 [36] database with 3,165 ICAO-compliant samples. In addition, 529 subjects with 1413 samples from the FERET [37] database of ICAO-compliant quality were used for the evaluation. From the samples, fea-

tures are extracted with the open-source feature extraction model ArcFace [12] which produces face templates of 512-dimensional floating point vectors with documented good performance on the used data set [4]. For the stable hash generation using k-means clustering, the parameters $P = 1$ subspace and $K = 64$ clusters are chosen in accordance with the size of the database. The experiments were implemented in Python and C++ on macOS Monterey 12.4 with an M2 processor at 3.50 GHz CPU clock frequency.

For the homomorphic operations, the CKKS [9] FHE scheme was applied, as it does not impair the biometric performance. The implementation of the state-of-the-art FHE library OpenFHE [2] was applied, where CKKS parameters corresponding to 128 bits of security were chosen [1]. For further workload reduction, coefficient packing for a quadratic speed-up as previously proposed by [4] was applied, showing the compatibility of HEBI with such approaches. The squared Euclidean distance was applied as the comparison metric. For the lattice-based PEKS scheme, the implementation by [5] was used.

6.1. Results

The results of the experimental evaluation are presented in Tables 1 and 2. In terms of execution times (Table 1), it can be seen that the majority of the workload is absorbed by the FHE comparisons on the encrypted feature vectors, an observation which is consistent with related work [16, 17, 19]. It is important to note that this workload can differ for different FHE schemes and has generally been found to be lower for integer-quantised and binary encryption, which introduces a trade-off with the biometric performance [30]. The baseline and preselection accuracy can be seen in Table 2, where a closed-set identification scenario was evaluated. Aside from this concrete instantiation however, we stress that HEBI is independent of the concrete preselection procedure and inherits and maintains the accuracy of the underlying indexing algorithm in question.

The main focus of this evaluation is the overhead of a

Table 1. HEBI execution times for 533 subjects and 64 clusters.

System function	Time (ms)
Probe stable hash generation	0.28
Probe encryption	2.27
PEKS search	7.69
FHE comparisons	9,996.00
Total	10,006.24
Baseline (exh. search)	334,891.00

Table 2. Accuracy of the stable hash clustering [34] for the FERET [37] and FRGCv2 [36] databases and $K = 64$ clusters.

Database	Enroll Samples	Search Samples	False Negative	True Positive	Presel. Accuracy	Baseline Accuracy
FERET [37]	529	884	19	865	0.9785	1.0000
FRGCv2 [36]	533	2,632	207	2,425	0.9214	0.9971

secure indexing using HEBI over unprotected preselection. From Table 1, it can be derived that the protected preselection using PEKS takes 7.69 milliseconds for 64 clusters or 0.12 milliseconds per cluster. As the cost for the preselection scales linearly with the number of clusters rather than the size of the enrolment database, this cost is expected to grow significantly slower than the cost for an exhaustive identification search. For larger databases, the original work on stable hashing [34] proposes a number of $K = 1024$ clusters, the cost of which can be approximated at 123.04 milliseconds, which can be considered real-time. Depending on the indexing algorithm used, there exists a trade-off between the number of clusters, the preselection error, and the number of biometric references per cluster. Overall, it becomes evident however that the lattice-based PEKS scheme adds only a negligible overhead to the identification system at less than 8% of the total cost, while providing post-quantum protection under preselection. Compared to the baseline system, the workload is reduced down to 3%. The communication cost for HEBI consists of 2.66MB for a CKKS public key, 267.4KB for a CKKS ciphertexts, 27.2KB for a PEKS public key, 52KB for a PEKS ciphertext, and 27KB for a PEKS trapdoor.

6.2. Security Analysis

The security of both the FHE and the PEKS scheme are based on the R-LWE [31] problem, which is assumed to be post-quantum secure. The HEBI protocol maintains the post-quantum security through all steps of the identification transaction, including preselection. Contrary to unprotected indexing approaches such as [13, 34, 38], the PEKS ciphertexts are generated in a non-deterministic manner, which makes them indistinguishable over the given clusters. A privacy attack as discussed in Section 3 is thereby prevented.

With regards to the requirements formulated in ISO/IEC 24745 [25], irreversibility is given through the security as-

sumption of R-LWE [31]. Unlinkability and renewability can be derived directly from the IND-CPA security of both the CKKS [9] and PEKS [5] schemes, i.e., the indistinguishability under chosen plaintext attacks. Through this property, an attacker cannot distinguish between an encryption of 0 and an encryption of 1. In biometric identification, this extends to the indistinguishability of encrypted templates: even if an attacker gains access to two encryptions of the same template, they cannot be distinguished from arbitrary inputs in a feasible manner. The same property holds for the encryption of index strings through the PEKS scheme. Therefore, it is not possible for an attacker to link data subjects to other subjects enrolled under the HEBI protocol or another system.

Finally, the performance preservation of HEBI is given through the application of CKKS [9] and PEKS [5], as neither scheme impairs the biometric performance. The operations in the encrypted domain correspond directly to the operations in an unprotected biometric system. In terms of computational performance of HEBI, our experimental evaluation has shown that the overhead of the PEKS scheme is small, while a trade-off between the preselection error and homomorphic workload persists. Further limitations of HEBI include the assumption of the semi-honest adversary model. Although this is an established assumption in biometric template protection, it does not fully reflect the capabilities of real-world adversaries. In addition, we have only evaluated the efficiency of HEBI for fixed-length feature representations, which can be considered a limitation.

7. Conclusion

This work firstly revealed that indexing schemes can leak privacy-sensitive biometric information. Motivated by this, we introduced the HEBI protocol for biometric indexing in the encrypted domain. Index strings in biometric identification systems allow for the reconstruction of privacy-sensitive information about the data subjects, which stands in violation to ISO/IEC 24745 as well as the GDPR. As a solution to this problem, HEBI gives post-quantum secure protection to the feature vectors alongside their auxiliary indexing data in storage, preselection, and comparison. HEBI is independent of the indexing algorithm and protection of the enrolment database and adds only negligible computational overhead per indexing cluster.

Acknowledgment

This research work has been funded by the German Federal Ministry of Education and Research and the Hessian Ministry of Higher Education, Research, Science and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE.

References

- [1] M. Albrecht, M. Chase, H. Chen, J. Ding, S. Goldwasser, S. Gorbunov, S. Halevi, J. Hoffstein, K. Laine, K. Lauter, S. Lokam, D. Micciancio, D. Moody, T. Morrison, A. Sahai, and V. Vaikuntanathan. Homomorphic encryption security standard. Technical report, HomomorphicEncryption.org, Toronto, Canada, November 2018.
- [2] A. A. Badawi, J. Bates, F. Bergamaschi, D. B. Cousins, S. Erabelli, N. Genise, S. Halevi, H. Hunt, A. Kim, Y. Lee, Z. Liu, D. Micciancio, I. Quah, Y. Polyakov, S. R.V., K. Rohloff, J. Saylor, D. Suponitsky, M. Triplett, V. Vaikuntanathan, and V. Zucca. OpenFHE: Open-source fully homomorphic encryption library. Cryptology ePrint Archive, Paper 2022/915, 2022.
- [3] P. Bauspieß, J. Kolberg, P. Drozdowski, C. Rathgeb, and C. Busch. Privacy-preserving preselection for protected biometric identification using public-key encryption with keyword search. *IEEE Transactions on Industrial Informatics*, 2022.
- [4] P. Bauspieß, J. Olafsson, J. Kolberg, P. Drozdowski, C. Rathgeb, and C. Busch. Improved homomorphically encrypted biometric identification using coefficient packing. In *Proc. Intl. Workshop on Biometrics and Forensics (IWBF)*, 2022.
- [5] R. Behnia, A. A. Yavuz, and M. O. Ozmen. High-speed high-security public key encryption with keyword search. In *IFIP Ann. Conf. on Data and Applications Security and Privacy*, pages 365–385. Springer, 2017.
- [6] D. Boneh, G. Di Crescenzo, R. Ostrovsky, and G. Persiano. Public key encryption with keyword search. In *Intl. Conf. on the Theory and Applications of Cryptographic Techniques*, pages 506–522. Springer, 2004.
- [7] D. Boneh and M. Franklin. Identity-based encryption from the Weil pairing. In *Proc. Annual Intl. Cryptology Conf.*, pages 213–229. Springer, 2001.
- [8] Z. Brakerski. Fully homomorphic encryption without modulus switching from classical GapSVP. In *Proc. Annual Intl. Cryptology Conf.*, pages 868–886. Springer, 2012.
- [9] J. H. Cheon, A. Kim, M. Kim, and Y. Song. Homomorphic encryption for arithmetic of approximate numbers. In *Intl. Conf. on the Theory and Appl. of Crypt. and Information Security*, pages 409–437. Springer, 2016.
- [10] I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène. TFHE: Fast fully homomorphic encryption over the torus. *Journal of Cryptology*, 33(1):34–91, 2020.
- [11] I. Csiszár. I-divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, pages 146–158, 1975.
- [12] J. Deng, J. Guo, and S. Zafeiriou. ArcFace: Additive angular margin loss for deep face recognition. In *Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [13] X. Dong, S. Kim, Z. Jin, J. Y. Hwang, S. Cho, and A. B. J. Teoh. Open-set face identification with index-of-max hashing by learning. *Pattern Recognition*, 103:107277, 2020.
- [14] P. Drozdowski, N. Buchmann, C. Rathgeb, M. Margraf, and C. Busch. On the application of homomorphic encryption to face identification. In *Intl. Conf. of the Biometrics Special Interest Group (BIOSIG)*, pages 173–180. Gesellschaft für Informatik e.V., September 2019.
- [15] P. Drozdowski, C. Rathgeb, and C. Busch. Computational workload in biometric identification systems: An overview. *IET Biometrics*, 8(6):351–368, November 2019.
- [16] P. Drozdowski, C. Rathgeb, and C. Busch. Turning a vulnerability into an asset: Accelerating facial identification with morphing. In *Intl. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 2582–2586. IEEE, May 2019.
- [17] P. Drozdowski, F. Stockhardt, C. Rathgeb, D. Osorio-Roig, and C. Busch. Feature fusion methods for indexing and retrieval of biometric data: Application to face recognition with privacy protection. *IEEE Access*, 9:139361–139378, October 2021.
- [18] L. Ducas, V. Lyubashevsky, and T. Prest. Efficient identity-based encryption over NTRU lattices. In *Intl. Conf. on the Theory and Application of Cryptology and Information Security*, pages 22–41. Springer, 2014.
- [19] J. J. Engelsma, A. K. Jain, and V. N. Boddeti. HERS: Homomorphically encrypted representation search. *IEEE Trans. on Biometrics, Behavior, and Identity Science (T-BIOM)*, 2022.
- [20] European Council. Regulation of the European Parliament and of the Council on the protection of individuals with regard to the processing of personal data and on the free movement of such data (General Data Protection Regulation), April 2016.
- [21] J. Fan and F. Vercauteren. Somewhat practical fully homomorphic encryption. Cryptology ePrint Archive, Paper 2012/144, 2012.
- [22] C. Gentry. Fully homomorphic encryption using ideal lattices. In *ACM Symposium on Theory of Computing*, pages 169–178, 2009.
- [23] M. Gomez-Barrero, E. Maiorana, J. Galbally, P. Campisi, and J. Fierrez. Multi-biometric template protection based on Homomorphic Encryption. *Pattern Recognition*, 67:149–163, July 2017.
- [24] J. Hämmerle-Uhl, G. Penn, G. Pötzelsberger, and A. Uhl. Size-reduction strategies for iris codes. *Intl. Journal of Computer and Information Engineering*, 9(1):290–293, 2015.
- [25] ISO/IEC JTC1 SC27 Security Techniques. *ISO/IEC 24745:2022. Information Technology - Security Techniques - Biometric Information Protection*. International Organization for Standardization, 2022.
- [26] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021.
- [27] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.
- [28] R. Kessler, O. Henninger, and C. Busch. Fingerprints, forever young? In *Proc. Intl. Conf. on Pattern Recognition (ICPR)*, pages 8647–8654, 2021.
- [29] J. Kolberg, P. Bauspieß, M. Gomez-Barrero, C. Rathgeb, M. Dürmuth, and C. Busch. Template protection based on homomorphic encryption: Computationally efficient application to iris-biometric verification and identification.

- In *IEEE Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2019.
- [30] J. Kolberg, P. Drozdzowski, M. Gomez-Barrero, C. Rathgeb, and C. Busch. Efficiency analysis of post-quantum-secure face template protection schemes based on homomorphic encryption. In *Intl. Conf. of the Biometrics Special Interest Group (BIOSIG)*, pages 175–182. Gesellschaft für Informatik e.V., September 2020.
- [31] V. Lyubashevsky, C. Peikert, and O. Regev. On ideal lattices and learning with errors over rings. In *Proc. 29th Ann. Intl. Conf. on the Theory and Applications of Cryptographic Techniques*, pages 1–23. Springer, 2010.
- [32] J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proc. 5th Berkeley Symposium on Math., Stat., and Prob.*, page 281, 1965.
- [33] T. Murakami, R. Fujita, T. Ohki, Y. Kaga, M. Fujio, and K. Takahashi. Cancelable permutation-based indexing for secure and efficient biometric identification. *IEEE Access*, 7:45563–45582, 2019.
- [34] D. Osorio-Roig, C. Rathgeb, P. Drozdzowski, and C. Busch. Stable hash generation for efficient privacy-preserving face identification. *Trans. on Biometrics, Behavior, and Identity Science (TBIOM)*, 4(3):333–348, July 2021.
- [35] A. Pflug, C. Rathgeb, U. Scherhag, and C. Busch. Binarization of spectral histogram models: An application to efficient biometric identification. In *2015 IEEE 2nd International Conference on Cybernetics (CYBCONF)*, pages 501–506. IEEE, 2015.
- [36] J. Phillips, P. Flynn, T. Scruggs, K. Bowyer, J. Chang, et al. Overview of the Face Recognition Grand Challenge. In *Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 947–954. IEEE, June 2005.
- [37] J. Phillips, H. Moon, S. Rizvi, and P. Rauss. The FERET evaluation methodology for face-recognition algorithms. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 22(10):1090–1104, October 2000.
- [38] A. Sardar, S. Umer, C. Pero, and M. Nappi. A novel cancelable facehashing technique based on non-invertible transformation with encryption and decryption template. *IEEE Access*, 8:105263–105277, 2020.
- [39] W. J. Scheirer and T. E. Boult. Cracking fuzzy vaults and biometric encryption. In *2007 Biometrics Symposium*, pages 1–6. IEEE, 2007.
- [40] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 1948.
- [41] J. Surbiryala, R. Raghavendra, and C. Busch. Finger vein indexing based on binary features. In *2015 Colour and Visual Computing Symposium (CVCS)*, pages 1–6. IEEE, 2015.
- [42] Unique Identification Authority of India. Aadhaar Dashboard. https://www.uidai.gov.in/aadhaar_dashboard/. Accessed 2023-04-26.
- [43] Y. Wang, J. Wan, J. Guo, Y.-M. Cheung, and P. C. Yuen. Inference-based similarity search in randomized montgomery domains for privacy-preserving biometric identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(7):1611–1624, 2017.
- [44] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshihara. Packed homomorphic encryption based on ideal lattices and its application to biometrics. In *Intl. Conf. on Availability, Reliability, and Security*, pages 55–74. Springer, 2013.
- [45] X. Zhang, C. Huang, D. Gu, J. Zhang, and H. Wang. BIB-MKS: post-quantum secure biometric identity-based multi-keyword search over encrypted data in cloud storage systems. *IEEE Transactions on Services Computing*, 2021.
- [46] Y. Zhang, J. Qin, and L. Du. A secure biometric authentication based on PEKS. *Concurrency and Computation: Practice and Experience*, 28(4):1111–1123, 2016.