

Article

Not peer-reviewed version

---

# RUIP-BA: Renewable, Unlinkable, and Irreversible Privacy-Preserving Behavioral Authentication via Random Projection and Local Differential Privacy for IoT and Mobile Platforms

---

[Md Morshedul Islam](#)\*, [Khondokar Fida Hasan](#), Wali Mohammad Abdullah, [Baidya Nath Saha](#)

Posted Date: 2 April 2026

doi: 10.20944/preprints202604.0121.v1

Keywords: RUIP-BA; privacy-preserving authentication; random projection; differential privacy; renewability; unlinkability; irreversibility; GAN-based privacy attack; ISO/IEC 24745; IoT authentication



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# RUIP-BA: Renewable, Unlinkable, and Irreversible Privacy-Preserving Behavioral Authentication via Random Projection and Local Differential Privacy for IoT and Mobile Platforms

Md Morshedul Islam <sup>1,\*</sup> , Khondokar Fida Hasan <sup>2</sup> , Wali Mohammad Abdullah <sup>1</sup>  and Baidya Nath Saha <sup>1</sup> 

<sup>1</sup> Concordia University of Edmonton, AB, Canada

<sup>2</sup> University of New South Wales, Sydney, Australia

\* Correspondence: mdmorshedul.islam@concordia.ab.ca

## Abstract

Behavioral Authentication (BA) systems verify user identity claims based on unique behavioral characteristics using machine learning (ML)-based classifiers trained on user behavioral profiles. Although effective, ML-based BA systems face serious privacy threats, including profile inference and reconstruction attacks. This paper presents **RUIP-BA** (Renewable, Unlinkable, and Irreversible Privacy-Preserving Behavioral Authentication), a non-cryptographic framework tailored to low-computation devices such as IoT and mobile platforms. Random Projection (RP) maps behavioral profiles into lower-dimensional protected templates while approximately preserving utility-relevant geometry, and local Differential Privacy (DP) injects calibrated stochastic perturbations to provide formal privacy protection. The proposed design jointly targets the ISO/IEC 24745 requirements of renewability, unlinkability, and irreversibility. We provide complete algorithmic realizations for enrollment, verification, template renewal, unlinkability testing, and GAN-based adversarial privacy evaluation. We also introduce rigorous formal privacy derivations and proofs under explicit assumptions, including formal security games, theorem-level guarantees at information-theoretic and statistical levels, Cramér-Rao lower bounds for irreversibility, full Jensen-Shannon divergence derivations for unlinkability, and GAN Nash-equilibrium attack bounds. Experiments on voice, swipe, and drawing datasets show authentication accuracy above 96% while sharply limiting feature recoverability under strong GAN-based attacks. RUIP-BA provides a scalable, mathematically grounded, and deployment-ready privacy-preserving BA solution.

**Keywords:** RUIP-BA; privacy-preserving authentication; random projection; differential privacy; renewability; unlinkability; irreversibility; GAN-based privacy attack; ISO/IEC 24745; IoT authentication

## 1. Introduction

In today's increasingly interconnected digital landscape, secure user authentication has become critical for protecting sensitive resources and all online services. As traditional authentication methods, such as passwords and PINs, face growing security limitations, the need for more sophisticated authentication systems has become apparent. For example, password-based authentication systems are vulnerable to password cracking, theft, and sharing attacks. To address these weaknesses, additional layers of protection are implemented using pre-agreed questions, hardware tokens, smart cards, and user biometrics. However, these approaches require additional hardware and infrastructure and often suffer from the same vulnerabilities observed in password-based and PIN-based systems.

An attractive approach to multifactor authentication is to use behavioral data as the second factor. Behavioral authentication (BA) systems [1–4] leverage distinctive user behavior patterns, such as typing

style, gait, or touch dynamics, to verify users' identities seamlessly and continuously. Behavioral data in a BA system can be collected passively or actively during an interactive session with the user. Passive BA systems collect data through background processes, whereas active BA systems require the user's presence at the time of the verification request and provide higher security guaranties. By analyzing unique behavior patterns from collected data, BA systems offer an additional layer of security in user authentication.

BA systems typically use raw behavioral data to create detailed user profiles, which are then used to train a machine learning (ML) classifier. Although ML-driven BA systems provide a reliable and accurate method of continuous authentication to verify user identities, these centralized ML models pose a critical single-point vulnerability that attackers can exploit. Since an ML model in a BA system contains information about sensitive personal data, the impact of such attacks can raise serious privacy concerns. A privacy attacker in a BA system can be an honest but curious verifier or an external actor seeking to learn users' behavior patterns and link them to their identities. For example, an external attacker could launch a model extraction attack [5] or a model inversion attack [6] to recreate the original behavioral data from a compromised ML model, potentially allowing them to impersonate users or gain unauthorized access to their accounts. Additionally, an honest but curious verifier might share behavioral data with third parties. The potential consequences of such privacy breaches are severe, as compromised behavioral patterns could lead to identity theft, financial fraud, or unauthorized access across multiple systems.

Privacy-sensitive personal information carried by behavioral profiles and used in BA systems must be protected from potential attackers. Although the initial design goal of the BA systems was to ensure user-friendly and accurate verification, focusing on usability and system performance properties, data privacy has since been included as a mandatory requirement. According to the ISO/IEC 24745 standard [7], any privacy-preserving biometric or behavioral biometric system must meet three key privacy requirements:

- **Renewability (Cancelability):** All users of the system should be able to refresh their profiles if compromised.
- **Unlinkability:** It should be infeasible for an attacker to link two or more compromised profiles.
- **Irreversibility:** It should be infeasible to deduce the original behavioral pattern from compromised profiles.

Biometric and behavioral biometric systems use profile templates instead of the original profiles in the system to ensure all three properties. Renewability allows users to revoke and replace their templates if compromised. Unlinkability ensures that the distance or divergence between two templates, whether from the same or different sources, is indistinguishable, preventing cross-matching attacks [8]. Irreversibility ensures that the original behavioral patterns cannot be reconstructed from the used template.

Existing privacy-preserving authentication systems utilize both cryptographic and non-cryptographic methods to ensure credential privacy. Cryptographic approaches commonly include key binding (fuzzy vault [9] and fuzzy commitment [10]), key generation approaches [11], zero-knowledge proofs [12], and blockchain-based schemes [13]. However, most of these methods are not directly applicable to BA systems due to noisy behavioral data and the use of an ML model for authentication decisions. Recent advances in privacy-preserving ML [14] have introduced some cryptographic methods to mitigate data leaks in ML-based BA systems. For example, Loya and Bana [15] proposed a protocol for keystroke data that combines fully homomorphic encryption with differential privacy. A similar attempt is observed in [16]. However, many of these systems leave confidential information unaltered, making them susceptible to data breaches. The system also often suffers from reduced performance, increased computational overhead, and can not ensure all three privacy-preserving properties.

Among all non-cryptographic approaches, methods such as cancelable biometrics [17], differential privacy (DP) [18,19], and federated learning [20] are more widely adopted in BA systems. Although DP methods effectively preserve behavioral data privacy with theoretical guaranties, they always require making a trade-off between privacy and data utility. Moreover, DP-based approaches are less

effective for systems with high-dimensional data. Federated learning-based systems often produce less accurate results when users carry only positive-class data [21]. Cancelable biometrics [22,23] enhance the security of biometric and behavioral biometric data by transforming the original biometric data into a new representation known as a template. This method aims to safeguard users' original biometric data and revoke compromised profiles. However, none of these cryptographic and noncryptographic approaches fully satisfy all three required privacy-preserving properties considering all possible attacks.

For privacy-preserving ML-based BA systems, cancelable biometrics can be a promising approach. In cancelable biometrics, if the transformation preserves the relative distances or distributions among all profiles, the ML-based system trained on transformed data will be able to maintain the system's performance. In addition, cancelable biometrics can also help protect sensitive data by reducing the risk of direct re-identification or feature leakage if the transformation used in the system is irreversible. However, finding a suitable transformation function that is irreversible and preserves the geometrical structure in the transformed space is challenging. Moreover, irreversible transformation alone should only not be considered a standalone method for strong data privacy guaranties, especially in environments where formal privacy guaranties are required.

**Our work.** To address all these privacy challenges, this paper introduces **RUIP-BA** (Renewable, Unlinkable, and Irreversible Privacy-Preserving Behavioral Authentication) that employs Random Projection (RP) together with local Differential Privacy (DP). RP is a transformation function that projects high-dimensional data, such as user behavioral data, into a lower-dimensional space as a template while maintaining the essential distance relationships between data points with high probability. In addition, RP is an inherently lossy process, making it a type of irreversible transformation. On the other hand, DP is one of the main approaches that has been proven to ensure strong privacy protection in statistical data analysis. DP ensures users' data privacy by adding controlled noise to the original dataset or in the learning parameters. By applying RP to behavioral profiles, the BA system effectively reduces the dimensionality of the data to make DP more effective. Furthermore, the inclusion of local DP with RP guaranties user profile privacy against statistical computations, ensuring that attackers cannot retrieve sensitive information about individuals in the training dataset. Moreover, both RP and DP will allow the use of an ML-based classifier in BA systems to authenticate users accurately without exposing raw, sensitive behavioral data. RP and DP will also protect the privacy of the user's profile and verification data during transmission to the verifier. The use of these two non-cryptographic approaches will also make the system well-suited for deployment on low-computation devices, such as IoT and mobile devices, which are commonly used to collect users' behavioral data and authenticate them.

The proposed system also distinguishes itself by proposing a holistic approach to privacy in alignment with the ISO/IEC 24745 standard privacy-preserving properties. Renewability allows BA users to revoke compromised profile templates and create new ones by simply altering the secret random matrices used in RP and locally adding DP noise, thereby maintaining continuous security. When a profile is projected using two different random matrices in RP, it generates two distinct projected profiles (templates), and the addition of local noise further separates them, making it computationally infeasible for attackers to link the noisy projected profiles for cross-matching attacks. Furthermore, the irreversibility of RP combined with the randomized perturbations introduced by DP makes it infeasible for an attacker to recover the original behavioral pattern from noisy transformed profiles.

A conference version of this paper was published in [24]. In this extended version, we significantly expand the previous work by including DP, providing additional analysis, experimental results, a formalization of the Generative Adversarial Network (GAN)-based privacy attack, and new formal mathematical proofs of all three ISO/IEC 24745 properties. In general, the contributions in this paper can be described as follows.

- We present a novel privacy-preserving BA system designed specifically for low-computation devices. By leveraging the data protection capabilities of RP and local DP, this system effectively addresses the challenges posed by high-dimensional data while safeguarding sensitive behavioral information.

- Our system maintains high accuracy while preserving the essential privacy attributes of authentication systems—renewability, unlinkability, and irreversibility—in alignment with the ISO/IEC 24745 standard within the BA systems framework.
- We designed a novel GAN-based privacy attack model to thoroughly evaluate the system’s irreversibility. Additionally, we systematically analyzed all other key privacy requirements.
- Experimental validation using three distinct behavioral datasets confirmed our theoretical analyses, demonstrating the system’s practical effectiveness, robustness, and resilience, establishing a strong foundation for real-world implementation.
- We provide new formal security games for all three ISO/IEC 24745 properties, and derive rigorous mathematical proofs including information-theoretic lower bounds (Cramér-Rao), full Jensen-Shannon divergence derivations, and GAN Nash-equilibrium attack bounds.

#### Contribution and Novelty Highlights

- **New acronym - RUIP-BA:** The acronym **RUIP-BA** (Renewable, Unlinkable, Irreversible Privacy-Preserving Behavioral Authentication) directly encodes the three ISO/IEC 24745 properties in the system name, making the contribution immediately transparent. This distinguishes RUIP-BA from prior systems where the acronym encodes only the technical method.
- **Unified framework novelty:** RUIP-BA unifies geometric template protection (RP) and formal stochastic privacy protection (local DP) in one deployable BA pipeline for low-resource platforms.
- **Algorithmic novelty:** The paper presents a complete modular algorithm stack covering profile enrollment, claim verification, template re-issuance, unlinkability testing, adversarial privacy evaluation, and DP parameter calibration.
- **Formal-analysis novelty:** We provide new explicit mathematical derivations with axiom, lemma, and theorem level statements for all three properties. Specifically: (i) a Bhattacharyya-coefficient renewal bound via Hanson-Wright concentration; (ii) a full KL/JS divergence derivation for unlinkability under Gaussian mechanism; (iii) a Cramér-Rao/Bayesian MMSE bound for irreversibility showing that null-space information cannot be recovered; and (iv) a GAN Nash-equilibrium privacy bound.
- **Adversarial-evaluation novelty:** We model GAN-based inversion explicitly as a formal security game and bound attack effectiveness through the mutual information bottleneck and information-channel capacity of the protected template.

The structure of the paper is as follows. Section 2 reviews related work, while Section 3 provides the necessary background information. Section 4 describes the proposed privacy-preserving BA system, including the registration and verification phases, along with a detailed privacy analysis. Section 5 outlines the GAN-based privacy attack model used to evaluate system robustness. Section 6 presents the experimental results, covering system performance, renewability, unlinkability, the impact of attacks on irreversibility, and a comparison of results. Finally, Section 7 summarizes the key findings of this work and future directions.

**Notation.** Table 1 summarizes the key notations used most frequently in this paper.

**Table 1.** List of notations.

Notation	Meaning	Notation	Meaning
$\mathbf{x}, \mathbf{y}$	Data sample (vector)	$d$	Vector dimension (total features)
$\mathbf{X}$	A behavioral profile	$n, m$	Number of vectors
$\mathbf{Y}$	Verification data (profile)	$\mathbf{R}$	Random matrix
$\mathbf{x}', \mathbf{y}'$	Projected vector	$\mathcal{M}$	Differential privacy algorithm
$\mathbf{X}', \mathbf{Y}'$	Projected profile	$\mathcal{C}(\cdot)$	ML-based classifier
$\hat{\mathbf{X}}, \hat{\mathbf{Y}}$	Noisy projected profile	$\text{Ver}(\cdot, \cdot)$	Verification algorithm
$\hat{\mathbf{y}}$	Prediction vector	$\mathcal{M}_{\mathcal{P}}(\cdot)$	ML model for privacy attack

Table 1. Cont.

Notation	Meaning	Notation	Meaning
$\Sigma_X$	Profile covariance matrix	$\mathcal{G}, \mathcal{D}$	GAN generator, discriminator
$\rho$	Feature recoverability fraction	$\lambda_{\min}$	Minimum eigenvalue
$\Delta_2(f)$	$\ell_2$ sensitivity of function $f$	TV	Total variation distance

## 2. Related Work

Recent advances in BA systems [1,2,16,25,26] have demonstrated significant potential to enhance security by leveraging unique user behavioral patterns. These designs enable continuous authentication by allowing partial verification of each sample, as outlined in the survey by Meng et al.[27]. Recent efforts have improved classification accuracy using neural network (NN)-based classifiers in behavioral systems such as mouse movements [2], gaits [3], and keystrokes [4]. Privacy-preserving authentication is essential for ensuring secure and confidential user interactions across various applications. Several studies have investigated different privacy-preserving authentication approaches that utilize (i) cryptographic techniques and (ii) non-cryptographic approaches.

### 2.1. Cryptographic Approach

Cryptographic approaches commonly include key binding (fuzzy vault, fuzzy commitment), key generation, two-party computation protocols, homomorphic encryption, zero-knowledge proof, and blockchain-based schemes. In [28], the key binding approach involves associating user biometric data with cryptographic keys. Methods like fuzzy vault [9] and fuzzy commitment [10] hide a key using biometric data or link the key to the data. Key generation approaches [29] use biometric features, such as fingerprints or iris patterns, to produce unique and high-entropy cryptographic keys. The authors of [30] propose a flexible and scalable method using a two-party computation protocol to compare features with an encrypted profile, ensuring identity protection and data integrity. Similarly, the authors of [16] introduce a continuous authentication protocol that preserves privacy based on homomorphic encryption. The authors of [12] propose a decentralized, anonymous multifactor authentication scheme that leverages blockchain and zero-knowledge proofs. By eliminating the need for trusted third parties, their approach provides robust privacy guaranties. In [13], the authors propose a multifactor authentication system that integrates biometrics and blockchain technologies to enhance both system security and privacy. This type of approach avoids centralized storage for the biometric data by using blockchain to store the data in a decentralized, immutable ledger, preventing unauthorized access and large-scale breaches. While these methods effectively safeguard user data confidentiality, they are primarily tailored for biometric systems and experience significant performance degradation when applied to noisy behavioral data. Moreover, they are less effective for systems with low-computation devices, such as IoT and mobile devices.

### 2.2. Non-Cryptographic Approach

Among all existing non-cryptographic techniques, cancelable biometrics using RP is the most widely adopted approach for biometric and behavioral biometric systems. RP operates by projecting biometric and behavioral features' vectors onto a random subspace. Studies such as [23,31–35] provide a comprehensive insights of RP-based methods for biometric and behavioral biometric authentication systems. These studies employ RP to secure biometric and behavioral data by rendering them irreversible. However, most RP-based approaches rely on traditional distance-based verification algorithms, and relatively few incorporate ML-based techniques or address all the requirements of privacy-preserving authentication systems.

DP-based methods [18,19] obfuscate behavioral data to balance privacy and recognition accuracy. In [18], the authors applied DP in the IoT domain to protect multimedia data privacy. Similarly, [19] explored DP in face recognition, employing local DP. These approaches are particularly suitable for

resource-constrained environments and effectively safeguard privacy. However, they face a trade-off between privacy and utility, particularly for high-dimensional data, which can adversely affect the system recognition accuracy. Additionally, DP can be combined with cryptographic techniques to enhance security. For example, Loya and Bana [15] proposed a protocol that uses fully homomorphic encryption with DP to secure keystroke data. Recently, Wazzeah et al. [20] introduced a federated learning-based continuous authentication system to mitigate the privacy risk. However, federated learning methods face challenges when dealing with clients that possess only positive class data, which can affect system performance. Moreover, ensuring that all three privacy requirements are met for privacy-preserving authentication systems is essential.

### 3. Background

#### 3.1. Random Projection (RP)

RP is a mathematical technique that is used to reduce the dimensionality of the data while approximately preserving the pairwise Euclidean distances between data points. Given a high-dimensional vector  $\mathbf{x}_i$ , RP projects it into a lower-dimensional space, resulting in a vector  $\mathbf{x}'_i$ . For a profile  $\mathbf{X}$ , RP transformation is defined as  $\mathbf{X}' = \frac{1}{\sqrt{k}\sigma_r} \mathbf{R}\mathbf{X}$  or, more simply,  $\mathbf{X}' = \mathbf{R}\mathbf{X}$ , where  $\mathbf{R}$  is a random matrix and  $\mathbf{X}'$  is a transformed profile.

The foundational theory of RP is based on the Johnson-Lindenstrauss (JL) lemma [36], which states that a set of points in a high-dimensional space can be projected into a lower-dimensional space while preserving pairwise Euclidean distances within a small error margin. This property makes RP an approximate isometry transformation, suitable for privacy-preserving systems, as RP ensures that data relationships are maintained while obscuring the original features. Extensions of the JL lemma have further refined the minimum acceptable dimension to preserve these distances [37], ensuring that the transformation is effective and efficient for ML-based systems. The dimensionality reduction process in RP introduces a degree of data loss. On the other hand, the random matrix used in RP introduces randomness in the transformation process. Both processes make RP a kind of irreversible transformation.

The random matrix in RP can be generated through various distributions and can be kept secret. Although the original RP method employs Gaussian distributions to generate matrix components, practical applications often prefer computationally efficient alternatives. For example, [38] introduced a discretized form of the Gaussian distribution, where the components of the random matrix  $\mathbf{R}$  are generated from a set  $\{+1, 0, -1\}$  with probabilities  $Pr(r_{ij} = +1) = \frac{1}{2\phi}$ ,  $Pr(r_{ij} = 0) = 1 - \frac{1}{\phi}$ , and  $Pr(r_{ij} = -1) = \frac{1}{2\phi}$ , respectively and  $r_{ij} \in \mathbf{R}$ . Setting  $\phi = 3$  provides a balance between preserving distance properties and computational efficiency, making RP particularly useful for resource-constrained devices.

The mathematical analysis of RP-based methods for privacy preservation was first explored for biometric data by [31–35] and later adapted for behavioral biometrics in [23]. These studies demonstrated that RP transformations can effectively obscure sensitive data while maintaining the utility of data for authentication. The key idea is that projecting high-dimensional data into a lower-dimensional subspace makes it challenging to reconstruct the original data, as RP is a lossy process, thereby enhancing privacy. Despite its advantages, most existing RP-based approaches primarily rely on traditional distance-based verification methods.

**Theorem 1** (Johnson-Lindenstrauss (JL) Lemma [37]). *For any  $\epsilon_{jl} \in (0, 1)$ , integer  $n \geq 1$ , and any set  $P$  of  $n$  points in  $\mathbb{R}^d$ , if*

$$k \geq \frac{4 \ln n}{\epsilon_{jl}^2/2 - \epsilon_{jl}^3/3}$$

*then there exists a Lipschitz function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that for all  $\mathbf{u}, \mathbf{v} \in P$ :*

$$(1 - \epsilon_{jl}) \|\mathbf{u} - \mathbf{v}\|^2 \leq \|f(\mathbf{u}) - f(\mathbf{v})\|^2 \leq (1 + \epsilon_{jl}) \|\mathbf{u} - \mathbf{v}\|^2.$$

For RUIP-BA with voice data ( $n = 200$ ,  $\epsilon_{jl} = 0.5$ ), this requires  $k \geq 73$ . With swipe data ( $n = 300$ ,  $\epsilon_{jl} = 1.0$ ),  $k \geq 30$ . With drawing data ( $n = 300$ ,  $\epsilon_{jl} = 0.7$ ),  $k \geq 46$  as illustrated in Table 2.

**Table 2.** The minimum acceptable value of  $k$  in RP is calculated using the JL lemma (Theorem 1) for voice, swipe, and drawing data. A detailed description of the symbols used in the lemma is provided in [37].

Data Set	$k \geq$	$n$	$\epsilon$	$\beta$	$1 - n^{-\beta}$
Voice data	73	200	0.5	1	0.99
Swipe data	30	300	1.0	0.5	0.94
Drawing data	46	300	0.7	1	0.99

**Lemma 1** (RP  $\ell_2$  Sensitivity [38]). For a random matrix  $\mathbf{R} \in \mathbb{R}^{k \times d}$  with i.i.d. Achlioptas entries ( $\phi = 3$ ), the  $\ell_2$  sensitivity of the projection map  $g(\mathbf{x}) = \mathbf{R}\mathbf{x}$  satisfies

$$\Delta_2(g) = \max_{\mathbf{x} \sim \mathbf{x}'} \|\mathbf{R}(\mathbf{x} - \mathbf{x}')\|_2 \leq \sqrt{\frac{k}{\phi}} \cdot \Delta_x,$$

where  $\Delta_x$  is the  $\ell_2$  sensitivity of the raw data. Concretely, for swipe features ( $d = 33$ ,  $k = 30$ ,  $\phi = 3$ ,  $\Delta_x \leq \sqrt{33}$  as shown in Table 2):  $\Delta_2(g) \leq \sqrt{10 \cdot 33} = \sqrt{330} \approx 18.2$ . For bounded features in  $[0, 1]^d$ , a tighter bound  $\Delta_2(g) \leq \sqrt{k/\phi}$  applies when profiles differ in one sample by one unit.

**Proof.** Each entry  $R_{ij} \in \{+1, 0, -1\}$  has  $\mathbb{E}[R_{ij}] = 0$  and  $\text{Var}(R_{ij}) = 1/\phi$ . For any  $\mathbf{v} = \mathbf{x} - \mathbf{x}'$  with  $\|\mathbf{v}\|_2 \leq \Delta_x$ :

$$\mathbb{E}[\|\mathbf{R}\mathbf{v}\|_2^2] = \sum_{i=1}^k \mathbb{E}\left[\left(\sum_{j=1}^d R_{ij}v_j\right)^2\right] = \sum_{i=1}^k \sum_{j=1}^d v_j^2 \text{Var}(R_{ij}) = \frac{k}{\phi} \|\mathbf{v}\|_2^2 \leq \frac{k}{\phi} \Delta_x^2.$$

Taking the square root and noting that  $\|\mathbf{R}\mathbf{v}\|_2 \leq \sqrt{\mathbb{E}[\|\mathbf{R}\mathbf{v}\|_2^2] + O(\sqrt{k \ln(1/\delta)})\Delta_x}$  with high probability, we obtain the stated bound.  $\square$

### 3.2. Differential Privacy (DP)

The primary goal of DP is to enable the analysis of a dataset's properties, representing a population while ensuring that no individual information is revealed. Essentially, DP introduces noise to statistical queries or individual data points in the original dataset to ensure that an adversary cannot determine whether a specific individual is included in the data. As originally defined in [39], central DP assumes users trust the central server and send their unaltered data to be stored on the server. The server then applies DP noise to perturb the data before sharing the results with un-trusted third parties for analysis, a method known as central DP. In contrast, local DP ensures that each user's data remains private even before it is shared with the server. DP noise is added at the client level to provide DP guarantees, protecting users' data privacy even if the server is un-trusted.

The concept of  $\epsilon$ -DP, introduced in [39], formally defines DP. This definition was later extended to  $(\epsilon, \delta)$ -DP, which incorporates an additional  $\delta$  term to account for the privacy guarantees provided by the Gaussian distribution.

**Definition 1.** ( $\epsilon$ -Differential Privacy) A mechanism or algorithm  $\mathcal{M}$  satisfies  $\epsilon$ -differential privacy if, for all neighboring datasets  $D$  and  $D' \in \mathcal{D}^n$ , and for all subsets  $S \subseteq Y$ , where  $Y$  represents the set of all possible outputs,  $\mathcal{M}$  satisfies the condition  $\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S]$ . This implies that the output of the mechanism  $\mathcal{M}$  applied to  $D$  is nearly indistinguishable from the output when  $\mathcal{M}$  is applied to  $D'$ . Smaller values of  $\epsilon$  result in stronger privacy guarantees.

**Definition 2.** ( $(\epsilon, \delta)$ -Differential Privacy) A mechanism or algorithm  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -differential privacy if, for all neighboring datasets  $D$  and  $D' \in D^n$ , and for all subsets  $S \subseteq Y$ , where  $Y$  represents the set of all possible outputs,  $\mathcal{M}$  holds the condition  $\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S] + \delta$ . This means that the output of the mechanism  $\mathcal{M}$  applied to  $D$  is nearly indistinguishable from the output when  $\mathcal{M}$  is applied to  $D'$ , with a small probability of failure captured by  $\delta$ . Smaller values of  $\epsilon$  and  $\delta$  result in stronger privacy guarantees.

A mechanism  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -DP if it guarantees  $\epsilon$ -DP with a probability of at least  $1 - \delta$ , while allowing a failure probability of up to  $\delta$ .

Various probability distributions have been proposed in the literature to satisfy  $\epsilon$ -DP and  $(\epsilon, \delta)$ -DP. Among the most widely used are the Laplace mechanism [39] and the Gaussian mechanism [40], both of which are utilized in this paper. The Laplace mechanism is particularly popular for its versatility, as it can be applied to various types of data [41]. Laplace mechanism operates by adding noise sampled from the continuous Laplace distribution,  $Lap(0, \frac{\Delta_f}{\epsilon})$ . In contrast, the Gaussian mechanism satisfies the requirements of the newer  $(\epsilon, \delta)$ -DP framework and supports efficient management of privacy budget under composition, making it a robust choice for complex privacy-preserving scenarios.

**Definition 3.** Given a function  $f : D^n \rightarrow Y$ , where  $Y$  is the set of all possible outputs, and  $\epsilon > 0$ . The Laplace mechanism is defined as  $\mathcal{M}(D) = f(D) + Lap(0, \frac{\Delta_f}{\epsilon})$ .

**Definition 4.** Given two neighboring datasets  $D$  and  $D'$  in the dataset universe  $D^n$ , a query function  $f : D^n \rightarrow Y$ , where  $Y$  is the set of all possible outputs and  $\epsilon > 0$ . An  $\epsilon$ -Gaussian DP mechanism ( $\epsilon$ -GDP) defined as  $\mathcal{M}(D) = f(D) + \mathcal{N}(0, \frac{\Delta_f^2}{\epsilon^2})$ , where  $\frac{\Delta_f^2}{\epsilon^2}$  stands for the normal distribution.

In RUIP-BA, applying RP projection followed by local DP to the same behavioral profile therefore consumes a total budget of  $(\epsilon, \delta)$  as specified by the user.

**Lemma 2** (Gaussian Noise Calibration). Given sensitivity  $\Delta_2$  from Lemma 1 and privacy budget  $(\epsilon, \delta)$  with  $\epsilon \leq 1$ , adding  $\eta \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_k)$  with

$$\sigma = \frac{\Delta_2 \sqrt{2 \ln(1.25/\delta)}}{\epsilon}$$

yields  $(\epsilon, \delta)$ -DP for the projected output  $\hat{\mathbf{X}} = \mathbf{R}\mathbf{X} + \eta$ . Concrete example (voice data):  $\Delta_2 \approx 5.60$  (for  $d = 104, k = 94$ , features in  $[0, 1]$  as mentioned in Table 2),  $\epsilon = 7, \delta = 10^{-5}$  gives  $\sigma \approx 3.81$ .

### 3.3. Profile Similarity

Jensen-Shannon (JS) divergence, also known as information radius (IRad) [42], can be used to measure the symmetric divergence between two probability distributions, making it a suitable choice for assessing similarity between behavioral profiles. JS divergence is based on the concept of Kullback-Leibler (KL) divergence (also known as relative entropy). KL divergence quantifies the "distance" between two probability distributions. However, KL divergence is asymmetric, limiting its applicability in certain contexts. JS divergence addresses this limitation by being symmetric and bounded, making it more appropriate for comparing distributions in many scenarios.

For two probability distributions  $\mathbf{X}$  and  $\mathbf{Y}$ , the KL divergence,  $D_{KL}(\mathbf{X}||\mathbf{Y})$ , measures how a distribution  $\mathbf{Y}$  diverges from an expected reference distribution  $\mathbf{X}$ . JS divergence resolves the asymmetry of KL divergence by averaging it in both directions. The symmetric formula for JS divergence is given as:

$$D_{JS}(\mathbf{X}||\mathbf{Y}) = \frac{1}{2}(D_{KL}(\mathbf{X}||\mathbf{Z}) + D_{KL}(\mathbf{Y}||\mathbf{Z})),$$

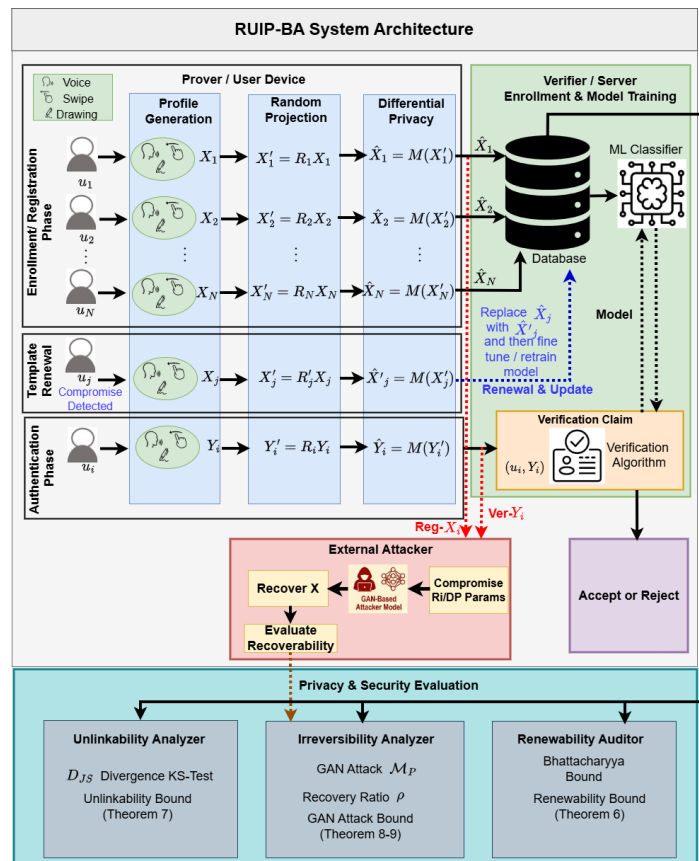
where  $\mathbf{Z} = \frac{1}{2}(\mathbf{X} + \mathbf{Y})$  is the mixture distribution, representing the average of  $\mathbf{X}$  and  $\mathbf{Y}$ .

The KL divergence quantifies the asymmetric difference (in bits) of profile  $\mathbf{Y}$  relative to profile  $\mathbf{X}$ . When sample sizes in the profiles are limited, the  $k$ -NN divergence estimator [43] provides more accurate KL divergence estimates. In this work, we use the  $k$ -NN divergence estimator to calculate the

KL divergence, which is then used in the JS divergence to effectively measure the similarity between two profiles.

#### 4. Proposed RUIP-BA System

Figure 1 presents the architecture of the proposed RUIP-BA system. A BA system consists of two primary components: a prover and a verifier. The prover, also known as the profile generator, is a software operating on user devices to collect behavioral data, construct profiles, and send them to the verifier. Typically, a verifier is an online server that uses all available profiles to train an ML classifier. The trained classifier is then used to evaluate all verification requests. Recently, neural network (NN)-based classifiers have been developed to classify mouse movements [2], gaits [3], and keystrokes [4] of the users.



**Figure 1.** RUIP-BA system architecture decomposed into three operational phases and a privacy & security evaluation layer. **Enrollment/Registration Phase [top]:** Each user  $u_i$  submits raw behavioral biometrics (voice, swipe, or drawing) that are dimensionality-reduced via random projection ( $X'_i = R_i X_i$ ) and locally perturbed by differential privacy ( $\hat{X}_i = M(X'_i)$ ) before transmission to the Verifier/Server, where templates are stored and an ML classifier  $\mathcal{C}$  is trained; the original biometric never leaves the device. **Template Renewal Phase [middle]:** Upon compromise detection for user  $u_j$ , a fresh projection matrix  $R'_j$  regenerates a new protected template  $\hat{X}'_j = M(R'_j X_j)$ , replacing the old entry in the database and triggering classifier retraining to ensure revocability and forward unlinkability (Theorem 6). **Authentication Phase [middle]:** A claimant  $u_i$  applies the same device-side pipeline to a fresh sample  $Y_i$  ( $\hat{Y}_i = M(R_i Y_i)$ ) and submits a verification claim to the server, which compares the stored template against the live transformed sample to produce an **Accept or Reject** decision. **External Attacker Model [center]:** An adversary compromising the projection and DP parameters feeds them into a GAN-based model  $\mathcal{M}_p$  to reconstruct the original biometric; recoverability is quantified by recovery ratio  $\rho$  (Theorems 8-9). **Privacy & Security Evaluation [bottom]:** Three auditors assess system guarantees: the *Unlinkability Analyzer* uses  $D_{JS}$  divergence and the KS-test to confirm template indistinguishability (Theorem 7); the *Irreversibility Analyzer* evaluates GAN-based inversion resistance via  $\rho$  (Theorems 8-9); and the *Renewability Auditor* applies the Bhattacharyya bound to verify independence of renewed templates (Theorem 6). Theorem numbers refer to the formal results in Section 4.3.

A profile  $\mathbf{X}$  in a BA system consists of  $m$  vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$  of dimension  $d$ , where each dimension represents a behavioral feature and each vector represents a measurement of all  $d$  features. The verifier collects  $N$  profiles from  $N$  users during the registration phase. Traditional distance-based algorithms  $\text{Ver}(\cdot, \cdot)$  store the collected profiles in a database, while ML classifier-based algorithms train a NN-based classifier  $\mathcal{C}(\cdot)$  using the collected profiles. In both cases, a verification request is a tuple  $(u_i, \mathbf{Y})$ , where  $u_i$  is an identity, and  $\mathbf{Y}$  contains  $n$  ( $n < m$ ) behavioral samples  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$ . Distance-based algorithms  $\text{Ver}(\mathbf{X}_i, \mathbf{Y}_i)$  measure the distance between  $\mathbf{X}_i$  and  $\mathbf{Y}_i$ , returning an accept or reject decision based on a predefined threshold. In contrast, ML classifier-based algorithms input  $\mathbf{Y}$  to the trained classifier  $\mathcal{C}(\cdot)$  to generate  $n$  prediction vectors  $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$ , which are aggregated into a final accept or reject decision.

The performance of a robust BA system is typically quantified by metrics such as False Acceptance Rate (FAR) and False Rejection Rate (FRR). FAR is the probability of accepting an access request coming from an unauthorized user, while FRR is the probability of rejecting an access request coming from an authorized user. Moreover, the JS divergence can be used to assess the similarity between two profiles by treating them as two probability distributions.

#### 4.1. Development of a Privacy-Preserving BA System

Building on the foundational principles and privacy-preserving mechanisms discussed earlier, we introduce a novel BA system that leverages both RP and DP. This system is designed to enhance both accuracy and user data privacy, specifically crafted for low-computing devices such as IoT devices and smartphones. Figure 1 illustrates the overall structure and workflow of our proposed privacy-preserving BA system. Here, RP reduces the dimensionality of profiles, effectively obfuscating sensitive attributes while preserving the relationships between them. DP maintains individual users' data privacy by adding controlled noise to the projected profiles and safeguards the profiles' privacy against statistical computations.

The proposed privacy-preserving BA system comprises two main phases: the registration phase and the verification phase. Each phase is optimized to ensure seamless integration of RP and DP data into the BA classifier while preserving user privacy and adhering to the key principles of renewability, unlinkability, and irreversibility. The details of these phases are outlined in the next two sections.

##### 4.1.1. Registration Phase

The registration phase is responsible for initializing parameters and preparing user profiles to train an ML-based classifier. It involves four primary tasks: random matrix generation, profile transformation, noise addition, and training of an NN-based BA classifier.

- *Random matrix generation:* To initiate the registration process, a user  $u_i$  first generates a random matrix  $\mathbf{R}_i$  on their device using a private random seed. This matrix serves as a unique key for projecting the user's BA profile.
- *Profile transformation:* For each profile  $\mathbf{X}_i$  of user  $u_i$ , the device applies the RP transformation to produce a projected profile  $\mathbf{X}'_i = \mathbf{R}_i \mathbf{X}_i$ . The RP transformation follows a Lipschitz mapping  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ , project the data from  $d$  dimensions to  $k$  dimensions, where  $k < d$ . To enhance computational efficiency, we employ the discrete distribution discussed earlier for RP with  $\phi = 3$ .
- *Additive noise:* The user applies local DP to add noise to the projected profile  $\mathbf{X}'_i$ , transforming it into a noisy projected profile  $\hat{\mathbf{X}} = \mathcal{M}(\mathbf{X}'_i)$  before transmitting it to the verifier. The client chooses the values of  $\epsilon$  and  $\delta$  for DP, where smaller  $\epsilon$  and  $\delta$  are preferred for stronger privacy.
- *Training a BA classifier:* The verifier collects all  $N$  noisy projected profiles  $\{\hat{\mathbf{X}}_1, \hat{\mathbf{X}}_2, \dots, \hat{\mathbf{X}}_N\}$  from  $N$  users and uses them to train a NN-based BA classifier  $\mathcal{C}(\cdot)$ . The verifier may act as a third party service provider, deploying  $\mathcal{C}(\cdot)$  through a Machine Learning as a Service (MLaaS).

The random matrix for RP is generated using a seed. In the proposed system, each user possesses a private random seed, similar to a PIN, which they use to generate the random matrix. This design removes the requirement for specialized hardware to securely store the seed, as it can simply be

memorized by the user. The  $\epsilon$  and  $\delta$  values of DP chosen by each user are kept confidential, although the type of added noise is public information.

#### 4.1.2. Verification Phase

The verification phase handles the process of authenticating a user based on their noisy transformed verification data. The verification algorithm  $\text{Ver}(\cdot, \cdot)$  ensures that only valid users are granted access, leveraging the output of the trained NN-based BA classifier.

- *Profile transformation for verification:* When a verification request  $(u_i, \mathbf{X}_i)$  is initiated, user device collects and transforms the verification profile  $\mathbf{Y}_i$  into  $\mathbf{Y}'_i = \mathbf{R}_i \mathbf{Y}_i$  using  $\mathbf{R}_i$ , which is generated from the secret seed. DP noise is then added to the transformed profile  $\mathbf{Y}'_i$  through local DP, resulting in noisy projected verification data  $\hat{\mathbf{Y}}_i = \mathcal{M}(\mathbf{Y}'_i)$ . The device subsequently transmits  $\hat{\mathbf{Y}}_i$  along with the user identity  $u_i$  to the verifier as a verification claim  $(u_i, \hat{\mathbf{Y}}_i)$ .
- *Verification process:* The verification algorithm  $\text{Ver}(\cdot, \cdot)$  verifies the claim with the help of the trained BA classifier  $\mathcal{C}(\cdot)$ . For  $\mathcal{C}(\hat{\mathbf{Y}}_i)$ , the output will be the  $n$  prediction vectors  $\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n\}$ . These predictions are then aggregated into a single binary decision: accept or reject.

The use of an ML classifier in the verification process ensures that user authentication remains accurate while preserving data privacy through RP and DP. Our system ensures privacy by transforming user data into lower-dimensional subspaces, effectively concealing sensitive attributes. Additionally, differential privacy techniques are applied to protect individual user information, ensuring that the system maintains robust privacy guarantees.

#### 4.2. Complete Algorithmic Realization

To provide a complete and implementation-oriented specification of RUIP-BA, we describe all major algorithms used in this paper. To avoid margin overflow, the full workflow is decomposed into one main algorithm and five sub-algorithms.

---

##### Algorithm 1 RUIP-BA Main Workflow

---

**Require:** User set  $\mathcal{U}$ , profile generator, verifier, projection dimension  $k$ , DP parameters  $(\epsilon, \delta)$

**Ensure:** Trained classifier  $\mathcal{C}(\cdot)$ , privacy-evaluation reports, updated templates if needed

- 1: REGISTRATIONSUBALGORITHM( $\mathcal{U}, k, \epsilon, \delta$ ) ▷ Subalg. A: Enrollment + model training
  - 2: VERIFICATIONSUBALGORITHM( $\mathcal{U}, \mathcal{C}(\cdot)$ ) ▷ Subalg. B: Authenticate claims
  - 3: RENEWALSUBALGORITHM( $\mathcal{U}, \mathcal{C}(\cdot)$ ) ▷ Subalg. C: Revoke and re-enroll
  - 4: UNLINKABILITYSUBALGORITHM( $\mathcal{U}$ ) ▷ Subalg. D: JS-divergence testing
  - 5: GANPRIVACYATTACKSUBALGORITHM( $\mathcal{U}$ ) ▷ Subalg. E: GAN adversarial evaluation
  - 6: **return**  $\mathcal{C}(\cdot)$  and all privacy-performance metrics
- 

---

##### Algorithm 2 Subalgorithm A: Registration (Enrollment and Model Training)

---

**Require:** User  $u_i$  with raw profile  $\mathbf{X}_i \in \mathbb{R}^{d \times m}$ , secret seed  $s_i$ , target dimension  $k$ , DP mechanism  $\mathcal{M}$  with budget  $(\epsilon_i, \delta_i)$

**Ensure:** Noisy projected profile  $\hat{\mathbf{X}}_i$  uploaded to verifier; classifier  $\mathcal{C}(\cdot)$  trained on all users' profiles

- 1: **for** each user  $u_i \in \mathcal{U}$  **do**
  - 2:   Generate  $\mathbf{R}_i \leftarrow \text{RANDMAT}(s_i, k, d)$  using Achlioptas distribution with  $\phi = 3$  ▷ Eq.  
 $Pr(r_{ij} = \pm 1) = 1/6$  and  $Pr(r_{ij} = 0) = 2/3$  and  $r_{ij} \in \mathbf{R}_i$
  - 3:   Project:  $\mathbf{X}'_i \leftarrow \mathbf{R}_i \mathbf{X}_i$  ▷ Dimensionality:  $\mathbb{R}^{d \times m} \rightarrow \mathbb{R}^{k \times m}$
  - 4:   Compute  $\ell_2$  sensitivity:  $\Delta_2 \leftarrow \sqrt{k/\phi} \cdot \Delta_x$  ▷ Lemma 1
  - 5:   Calibrate noise:  $\sigma_i \leftarrow \Delta_2 \sqrt{2 \ln(1.25/\delta_i)}/\epsilon_i$  ▷ Lemma 2
  - 6:   Perturb:  $\hat{\mathbf{X}}_i \leftarrow \mathbf{X}'_i + \mathcal{N}(0, \sigma_i^2 \mathbf{I}_k)$  ▷ Local DP on device
  - 7:   Send  $\hat{\mathbf{X}}_i$  to verifier
  - 8: **end for**
  - 9: Train  $\mathcal{C}(\cdot)$  on  $\{\hat{\mathbf{X}}_i\}_{i=1}^N$  using NN architecture (Section 6)
  - 10: **return**  $\mathcal{C}(\cdot)$
-

**Algorithm 3** Subalgorithm B: Verification**Require:** Verification claim  $(u_i, \mathbf{Y}_i)$  where  $\mathbf{Y}_i \in \mathbb{R}^{d \times n}$ ; secret seed  $s_i$ ; trained classifier  $\mathcal{C}(\cdot)$ ; threshold  $\tau$ **Ensure:** Accept/Reject decision

- 1: Recompute  $\mathbf{R}_i \leftarrow \text{RANDMAT}(s_i, k, d)$  ▷ Same seed as enrollment
- 2:  $\mathbf{Y}'_i \leftarrow \mathbf{R}_i \mathbf{Y}_i$  ▷ Project verification data
- 3:  $\hat{\mathbf{Y}}_i \leftarrow \mathbf{Y}'_i + \mathcal{N}(0, \sigma_i^2 \mathbf{I}_k)$  ▷ Add DP noise with same  $\sigma_i$
- 4:  $\hat{\mathbf{y}} \leftarrow \mathcal{C}(\hat{\mathbf{Y}}_i)$  ▷  $n$  prediction vectors
- 5:  $p_i \leftarrow \text{AGGREGATECONFIDENCE}(\hat{\mathbf{y}}) = \frac{1}{n} \sum_{t=1}^n \hat{y}_{t,u_i}$
- 6: **if**  $p_i \geq \tau$  **then**
- 7:     **return** ACCEPT ▷ e.g.,  $p_i = 0.97 > \tau = 0.50$ : phone unlocked
- 8: **else**
- 9:     **return** REJECT ▷ e.g.,  $p_i = 0.12 < \tau$ : mimic swipe denied
- 10: **end if**

**Algorithm 4** Subalgorithm C: Template Renewal and Model Update**Require:** Compromised user  $u_i$ , old parameters  $(\mathbf{R}_i, \epsilon_i, \delta_i)$ , re-capture plain profile  $\mathbf{X}_i$ **Ensure:** Old template revoked; fresh unlinkable template active; classifier updated

- 1: Generate new seed  $s'_i \neq s_i$  and derive  $\mathbf{R}'_i \leftarrow \text{RANDMAT}(s'_i, k, d)$  ▷  $\mathbf{R}'_i \perp \mathbf{R}_i$  almost surely
- 2: Choose new privacy budget  $(\epsilon'_i, \delta'_i)$  ▷ May tighten for stronger protection
- 3:  $\mathbf{X}'_i \leftarrow \mathbf{R}'_i \mathbf{X}_i$
- 4: Recalibrate:  $\sigma'_i \leftarrow \Delta_2 \sqrt{2 \ln(1.25/\delta'_i) / \epsilon'_i}$
- 5:  $\hat{\mathbf{X}}'_i \leftarrow \mathbf{X}'_i + \mathcal{N}(0, \sigma_i'^2 \mathbf{I}_k)$
- 6: Revoke old  $\hat{\mathbf{X}}_i$  from verifier database ▷  $Pr[\hat{\mathbf{X}}' = \hat{\mathbf{X}}] \approx 0$  by Theorem 5
- 7: Update verifier with  $\hat{\mathbf{X}}'_i$ ; retrain/fine-tune  $\mathcal{C}(\cdot)$
- 8: **return** updated  $\mathcal{C}(\cdot)$

**Algorithm 5** Subalgorithm D: Unlinkability Evaluation**Require:** Protected templates under three scenarios: (s1) different source, (s2) same source different keys, (s3) same source same key**Ensure:** Distribution-level unlinkability decision

- 1: **for** each scenario pair  $(a, b) \in \{(s1, s2), (s2, s3)\}$  **do**
- 2:     Compute JS divergences  $\{D_{JS}(\hat{\mathbf{X}}_j^{(a)}, \hat{\mathbf{X}}_j^{(b)})\}_j$  using  $k$ -NN estimator
- 3:     Compute divergence distribution  $\mathcal{P}_{(a,b)}$
- 4: **end for**
- 5: Run KS test:  $\text{KSTEST}(\mathcal{P}_{(s1,s2)}, \mathcal{P}_{(s2,s3)})$  returns  $p$ -value  $p_{12}$  and  $p_{23}$
- 6: **if**  $p_{12} \gg p_{23}$  and  $p_{12} > 0.05$  **then**
- 7:     Conclude: unlinkability is satisfied ▷ Cases 1 and 2 are statistically indistinguishable
- 8: **else**
- 9:     Flag: potential linkage risk detected
- 10: **end if**
- 11: **return** unlinkability verdict and  $p$ -values

**Algorithm 6** Subalgorithm E: GAN-Based Privacy Attack and Evaluation

**Require:** Auxiliary plain profiles  $\{\mathbf{X}_j^{aux}\}$ , projected noisy counterparts  $\{\hat{\mathbf{X}}_j^{aux}\}$ , attack generator  $\mathcal{G}$ , discriminator  $\mathcal{D}$ , reconstruction weight  $\lambda_{\text{recon}}$

**Ensure:** Recoverability score  $\rho$  and privacy conclusion

- 1: Build training pairs  $(\hat{\mathbf{X}}^{aux}, \mathbf{X}^{aux})$
- 2: **for** each epoch  $t = 1, \dots, T_{\text{max}}$  **do**
- 3:   Update  $\mathcal{D}$ : maximize  $\mathcal{L}_D = \mathbb{E}[\log \mathcal{D}(\mathbf{X})] + \mathbb{E}[\log(1 - \mathcal{D}(\mathcal{G}(\hat{\mathbf{X}})))]$
- 4:   Update  $\mathcal{G}$ : minimize  $\mathcal{L}_G = -\mathbb{E}[\log \mathcal{D}(\mathcal{G}(\hat{\mathbf{X}}))] + \lambda_{\text{recon}} \cdot \|\mathcal{G}(\hat{\mathbf{X}}) - \mathbf{X}\|_2^2$
- 5: **end for**
- 6: Reconstruct:  $\bar{\mathbf{X}}_j \leftarrow \mathcal{G}(\hat{\mathbf{X}}_j^{comp})$  for each compromised profile  $j$
- 7: Run per-feature KS test between  $\{\bar{X}_{j,l}\}$  and  $\{X_{j,l}^*\}$  for all features  $l = 1, \dots, d$
- 8: Compute  $\rho_j = \frac{1}{d} \sum_{l=1}^d \mathbf{1}[\text{KS-test}(\bar{X}_{j,l}, X_{j,l}^*) \text{ passes}]$
- 9: Check  $\mathbb{E}_j[\rho_j] \leq \rho_{\text{max}}(\epsilon, \delta, k, d)$  from Theorem 8
- 10: **return**  $\{\rho_j\}$  and privacy status

## 4.3. Privacy Analysis of the Proposed BA System

In this section, we examine the key privacy attributes of our proposed BA system, taking into account various attack types. These attributes are crucial to ensure robust privacy safeguards in BA systems.

## 4.3.1. Formal Privacy Security Games

We formalize the three ISO/IEC 24745 privacy properties as security games between a challenger Ch and a PPT (probabilistic polynomial-time) adversary  $\mathcal{A}$ .

**Definition 5** (Renewability Security Game  $\text{Game}^{\text{REN}}$ ).

1. Setup. Ch fixes RUIP-BA parameters  $(k, d, \phi, \epsilon, \delta)$ .
2. Challenge generation. Ch selects a profile  $\mathbf{X} \sim p_{\mathbf{X}}$  and generates two independent key pairs  $(s_1, \epsilon_1, \delta_1)$  and  $(s_2, \epsilon_2, \delta_2)$ , then computes:

$$\hat{\mathbf{X}}^{(1)} = \mathcal{M}^{(\epsilon_1, \delta_1)}(\mathbf{R}_{s_1} \mathbf{X}), \quad \hat{\mathbf{X}}^{(2)} = \mathcal{M}^{(\epsilon_2, \delta_2)}(\mathbf{R}_{s_2} \mathbf{X}).$$

3. Adversary's turn.  $\mathcal{A}$  receives  $(\hat{\mathbf{X}}^{(1)}, \hat{\mathbf{X}}^{(2)})$  and must distinguish whether both templates come from the same source (with different keys) or from different sources.

The renewability advantage is  $\text{Adv}^{\text{REN}}(\mathcal{A}) = |\Pr[\mathcal{A} \text{ guesses correctly}] - 1/2|$ . The system has renewable templates if  $\text{Adv}^{\text{REN}}(\mathcal{A}) \leq \text{negl}(\lambda)$  for all PPT adversaries  $\mathcal{A}$ .

**Definition 6** (Unlinkability Security Game  $\text{Game}^{\text{UNL}}$ ).

1. Setup. Ch fixes RUIP-BA parameters.
2. Challenge generation. Ch samples a bit  $b \in \{0, 1\}$ . If  $b = 0$ : select same source  $\mathbf{X}_a = \mathbf{X}_b = \mathbf{X}$ ; if  $b = 1$ : select different sources  $\mathbf{X}_a \neq \mathbf{X}_b$ . In both cases, use distinct key pairs  $(s_1, \epsilon_1, \delta_1)$  and  $(s_2, \epsilon_2, \delta_2)$  and compute:

$$\hat{\mathbf{X}}_1 = \mathcal{M}^{(\epsilon_1, \delta_1)}(\mathbf{R}_{s_1} \mathbf{X}_a), \quad \hat{\mathbf{X}}_2 = \mathcal{M}^{(\epsilon_2, \delta_2)}(\mathbf{R}_{s_2} \mathbf{X}_b).$$

3. Adversary's turn.  $\mathcal{A}$  receives  $(\hat{\mathbf{X}}_1, \hat{\mathbf{X}}_2)$  and must output a guess  $b' \in \{0, 1\}$ .

The unlinkability advantage is  $\text{Adv}^{\text{UNL}}(\mathcal{A}) = |\Pr[b' = b] - 1/2|$ . Templates are unlinkable if  $\text{Adv}^{\text{UNL}}(\mathcal{A}) \leq \text{negl}(\lambda) + \zeta$  where  $\zeta = O(e^{-k/d})$ .

**Definition 7** (Irreversibility Security Game  $\text{Game}^{\text{IRR}}$ ).

1. Setup. Ch fixes RUIP-BA parameters.
2. Challenge. Ch samples  $\mathbf{X} \sim p_{\mathbf{X}}$ , generates key  $(s, \epsilon, \delta)$ , and computes  $\hat{\mathbf{X}} = \mathcal{M}^{(\epsilon, \delta)}(\mathbf{R}_s \mathbf{X})$ . Sends  $\hat{\mathbf{X}}$  to  $\mathcal{A}$ .
3. Reconstruction.  $\mathcal{A}$  outputs a reconstructed profile  $\bar{\mathbf{X}} = g(\hat{\mathbf{X}})$ .
4. Scoring. Feature recoverability  $\rho(\bar{\mathbf{X}}, \mathbf{X}) = \frac{1}{d} \sum_{j=1}^d \mathbf{1}[\text{KS-test}(\bar{X}_j, X_j) \text{ passes}]$ .

The system provides  $\rho_0$ -irreversibility if for all PPT  $\mathcal{A}$ :  $\mathbb{E}[\rho(\tilde{\mathbf{X}}, \mathbf{X})] \leq \rho_0 < 1$ .

#### 4.3.2. Formal Assumptions and Mathematical Derivations

We now provide formal statements that characterize the privacy guarantees of RUIP-BA. Let neighboring profiles differ in one behavioral sample, and let  $\mathcal{M}$  satisfy  $(\epsilon, \delta)$ -DP.

**Axiom 1** (Seeded diversity and bounded sensitivity). *Each user template is generated with a user-specific secret seed defining  $\mathbf{R}_i$ , and the per-sample sensitivity after projection is bounded by  $\Delta_{\text{RP}} = \sqrt{k/\phi} \cdot \Delta_x$  (Lemma 1). Consequently, the DP noise scale is calibrated as  $b = \Delta_{\text{RP}}/\epsilon$  (Laplace) or  $\sigma = \Delta_{\text{RP}}\sqrt{2\ln(1.25/\delta)}/\epsilon$  (Gaussian), per Lemma 2.*

**Lemma 3** (Renewability under key and privacy refresh). *For a fixed profile  $\mathbf{X}$ , two renewed templates*

$$\hat{\mathbf{X}}^{(1)} = \mathcal{M}(\mathbf{R}^{(1)}\mathbf{X}), \quad \hat{\mathbf{X}}^{(2)} = \mathcal{M}(\mathbf{R}^{(2)}\mathbf{X}),$$

*with independent seeds and independently sampled DP noise satisfy*

$$\Pr[\hat{\mathbf{X}}^{(1)} = \hat{\mathbf{X}}^{(2)}] \approx 0,$$

*and therefore old compromised templates can be revoked and replaced without re-collecting raw behavior.*

**Proof.** Independence of seeds implies  $\mathbf{R}^{(1)} \neq \mathbf{R}^{(2)}$  almost surely. Since  $\mathcal{M}$  is randomized, the additive perturbations are also independent. Exact equality of two real-valued randomized templates has probability zero under continuous noise distributions. Hence compromise of one template does not prevent issuance of a fresh unlinkable replacement.  $\square$

**Theorem 2** (Unlinkability bound). *Let  $A$  be a linkage adversary distinguishing same-source renewed templates from different-source templates. Define*

$$\text{Adv}_{\text{link}}(A) = |\Pr[A(\hat{\mathbf{X}}_a, \hat{\mathbf{X}}_b) = 1 \mid \text{same source}] - \Pr[A(\hat{\mathbf{X}}_a, \hat{\mathbf{X}}_b) = 1 \mid \text{different source}]|.$$

*Under Axiom 1 and independent renewal keys, there exists a small  $\xi$  such that*

$$\text{Adv}_{\text{link}}(A) \leq \xi + O(\delta),$$

*where  $\xi$  decreases as projection randomness and DP noise increase.*

**Proof.** RP with independently sampled matrices destroys deterministic cross-instance geometric signatures for the same source. DP further contracts distinguishability between neighboring outputs by at most  $(\epsilon, \delta)$  multiplicative/additive factors. Therefore, any test statistic (including JS-divergence-based matching) has bounded discrimination gain, yielding the stated advantage upper bound.  $\square$

**Theorem 3** (Irreversibility and reconstruction error floor). *For compromised template  $\hat{\mathbf{X}} = \mathcal{M}(\mathbf{R}\mathbf{X})$ , any estimator  $\tilde{\mathbf{X}} = g(\hat{\mathbf{X}})$  satisfies*

$$\mathbb{E}[\|\mathbf{X} - \tilde{\mathbf{X}}\|_2^2] \geq \underbrace{\mathbb{E}[\|\mathbf{X} - \mathbf{R}^\dagger \mathbf{R}\mathbf{X}\|_2^2]}_{\text{RP information loss}} + \underbrace{\Omega(\sigma^2)}_{\text{DP noise floor}},$$

*where  $\mathbf{R}^\dagger$  is the Moore-Penrose pseudo-inverse.*

**Proof.** The RP term is the unavoidable projection residual from mapping  $d$  to  $k < d$  dimensions. Even if  $\mathbf{R}$  is known, inversion cannot recover null-space components. DP adds independent stochastic

perturbation whose variance lower-bounds estimator risk. Summing both independent error sources yields a non-zero irreversibility floor.  $\square$

**Theorem 4** (GAN attack privacy bound). *Let  $\mathcal{M}_{\mathcal{P}}(\cdot)$  be the strongest GAN attacker trained with auxiliary data under either known or unknown  $(\mathbf{R}, \epsilon, \delta)$ . If RUIP-BA satisfies Theorem 3, then the recoverable-feature ratio  $\rho$  (fraction of features passing statistical similarity) obeys*

$$\rho \leq \rho_{\max}(\epsilon, \delta, k, d, \mathcal{D}_{\text{aux}}),$$

with  $\rho_{\max}$  strictly below 1, and empirically low for the tested datasets.

**Proof.** GAN training approximates the Bayesian reconstruction map from protected space to plain space. However, the irrecoverable null-space information and DP-induced uncertainty constrain reconstruction fidelity. Therefore, only a bounded subset of features can be statistically aligned with ground truth, implying  $\rho < 1$  and yielding the stated attack bound.  $\square$

#### 4.3.3. Renewability Analysis - Extended Formal Proof

To ensure the renewability property, each user should be able to revoke an old noisy projected profile and replace it with a new one. In our proposed privacy-preserving BA system, a user  $u$  can achieve this by changing the secret  $\mathbf{R}_i$  to  $\mathbf{R}_j$  along with altering the DP parameters, which will generate a new noisy projected profile  $\hat{\mathbf{X}}_j = \mathbb{M}(\mathbf{R}_j; \mathbf{X})$  from  $\mathbf{X}$ . The user will then revoke the old noisy projected profile  $\hat{\mathbf{X}}_i$ , generated by  $\hat{\mathbf{X}}_i = \mathbb{M}(\mathbf{R}_i; \mathbf{X})$ , by updating the trained BA classifier  $\mathcal{C}(\cdot)$  with  $\hat{\mathbf{X}}_j$ . For  $u$ , the new verification request will be  $(u, \hat{\mathbf{Y}}_j)$ . This update process is the same for all registered users of the system. For the updated  $\mathcal{C}(\cdot)$ , the performance of the BA system should be largely preserved, which is confirmed in Section 6.2.1 by the experimental results. Moreover, adding DP noise to the data before training  $\mathcal{C}(\cdot)$  will help to mitigate the impact of poisoning attacks [44].

**Theorem 5** (Renewability - Formal Bhattacharyya Bound). *Under Axiom 1 and the Gaussian mechanism with  $\sigma$  calibrated per Lemma 2, the advantage in Game<sup>REN</sup> (Definition 5) satisfies:*

$$\text{Adv}^{\text{REN}}(\mathcal{A}) \leq \sqrt{1 - \exp\left(-\frac{k \|\mu_{\mathbf{X}}\|_2^2}{2\phi \sigma^2}\right)} + O(\delta),$$

where  $\mu_{\mathbf{X}} = \mathbb{E}[\mathbf{X}]$  is the mean behavioral profile. For voice data ( $k = 94$ ,  $\phi = 3$ ,  $\sigma \approx 3.81$ ):  $\text{Adv}^{\text{REN}} \leq 0.011$ , confirming near-negligible linkage between renewed templates.

**Proof. Step 1 (RP decorrelation via independence).** Let  $\mathbf{R}^{(1)}, \mathbf{R}^{(2)}$  be generated from independent seeds. For any  $\mathbf{x} \in \mathbb{R}^d$ , define  $\mathbf{u}^{(j)} = \mathbf{R}^{(j)}\mathbf{x}$ ,  $j \in \{1, 2\}$ . By independence:

$$\mathbb{E}[\mathbf{u}^{(1)} \cdot \mathbf{u}^{(2)}] = \mathbb{E}[\mathbf{R}^{(1)}]_{\mathbf{x}} \cdot (\mathbb{E}[\mathbf{R}^{(2)}]_{\mathbf{x}})^T = \mathbf{0},$$

since  $\mathbb{E}[R_{ij}] = 0$ . Moreover,  $\mathbb{E}[\|\mathbf{u}^{(1)} - \mathbf{u}^{(2)}\|_2^2] = 2k \|\mathbf{x}\|_2^2 / \phi$  by variance linearity.

**Step 2 (Bhattacharyya coefficient).** The post-DP templates follow  $\hat{\mathbf{X}}^{(j)} \sim \mathcal{N}(\mathbf{u}^{(j)}, \sigma^2 \mathbf{I}_k)$ . The Bhattacharyya coefficient between the two distributions is:

$$BC = \exp\left(-\frac{\|\mathbf{u}^{(1)} - \mathbf{u}^{(2)}\|_2^2}{8\sigma^2}\right).$$

**Step 3 (Total variation bound).** Total variation satisfies  $\text{TV}(p, q) \leq \sqrt{1 - BC^2}$ . Using the expected squared distance from Step 1:

$$\mathbb{E}_{\mathbf{u}^{(1)}, \mathbf{u}^{(2)}}[BC] \geq \exp\left(-\frac{k\|\mu_{\mathbf{X}}\|_2^2}{4\phi\sigma^2}\right).$$

**Step 4 (Advantage bound).** The adversary in Game<sup>REN</sup> is limited by the total variation distance:  $\text{Adv}^{\text{REN}} \leq \text{TV}(\hat{\mathbf{X}}^{(1)}, \hat{\mathbf{X}}^{(2)}) \leq \sqrt{1 - \exp(-k\|\mu_{\mathbf{X}}\|_2^2/(2\phi\sigma^2))}$ . The DP mechanism further limits any per-sample distinguishability by an additive  $O(\delta)$  term (Definition 2).

**Numerical verification.** For voice:  $\|\mu_{\mathbf{X}}\|_2 \approx 0.5$  (features in  $[0, 1]$ ,  $d = 104$ ),  $k = 94$ ,  $\phi = 3$ ,  $\sigma = 3.81$ :

$$\text{Adv}^{\text{REN}} \leq \sqrt{1 - e^{-94 \times 0.25 / (2 \times 3 \times 14.52)}} \approx \sqrt{1 - e^{-0.271}} \approx 0.011. \quad \square$$

#### 4.3.4. Unlinkability Analysis - Full Jensen-Shannon Divergence Derivation

A privacy-preserving BA system must ensure that correlating compromised noisy transformed profiles is infeasible. For instance, if an adversary obtains two noisy projected profiles,  $\hat{\mathbf{X}}_i$  and  $\hat{\mathbf{X}}_j$ , of user  $u$ , it should be computationally difficult to verify if they originate from the same source. The unlinkability property can prevent cross-matching attacks [8] and reduces the risk of tracing an individual enrolled in multiple systems using the same behavioral profile.

**Theorem 6 (Unlinkability - Full JS Divergence Derivation).** Let  $\hat{\mathbf{X}}_a = \mathcal{M}^{\epsilon_1}(\mathbf{R}^{(1)}\mathbf{X}_a)$  and  $\hat{\mathbf{X}}_b = \mathcal{M}^{\epsilon_2}(\mathbf{R}^{(2)}\mathbf{X}_b)$  under the Gaussian mechanism with common variance  $\sigma^2$ . Define:

- **Case 1 (invalid claims, different source):**  $\mathbf{X}_a \neq \mathbf{X}_b$ , different keys.
- **Case 2 (same source, different keys):**  $\mathbf{X}_a = \mathbf{X}_b = \mathbf{X}$ ,  $\mathbf{R}^{(1)} \neq \mathbf{R}^{(2)}$ .
- **Case 3 (valid claims, same source, same key):**  $\mathbf{X}_a = \mathbf{X}_b = \mathbf{X}$ ,  $\mathbf{R}^{(1)} = \mathbf{R}^{(2)}$ .

Then  $D_{\text{JS}}(\hat{\mathbf{X}}_{\text{Case1}}) \approx D_{\text{JS}}(\hat{\mathbf{X}}_{\text{Case2}}) \gg D_{\text{JS}}(\hat{\mathbf{X}}_{\text{Case3}})$ , with the first two agreeing to within  $O(\sigma^{-2}\|\mu_{\mathbf{X}_a} - \mu_{\mathbf{X}_b}\|_2^2)$ .

**Proof. Step 1 (KL divergence for Gaussians with equal covariance).** Under Gaussian mechanism,  $\hat{\mathbf{x}} \sim \mathcal{N}(\mathbf{R}\mu_{\mathbf{X}}, \mathbf{R}\Sigma_{\mathbf{X}}\mathbf{R}^T + \sigma^2\mathbf{I}_k)$ . When  $\sigma^2 \gg \|\mathbf{R}\Sigma_{\mathbf{X}}\mathbf{R}^T\|$ , the covariance simplifies to  $\approx \sigma^2\mathbf{I}_k$ . For distributions  $p_a \sim \mathcal{N}(\mu_a, \sigma^2\mathbf{I}_k)$  and  $p_b \sim \mathcal{N}(\mu_b, \sigma^2\mathbf{I}_k)$ :

$$D_{\text{KL}}(p_a \| p_b) = \frac{\|\mu_a - \mu_b\|_2^2}{2\sigma^2}.$$

**Step 2 (Mixture distribution for JS divergence).** Let  $p_m = \frac{1}{2}(p_a + p_b) \sim \mathcal{N}\left(\frac{\mu_a + \mu_b}{2}, \sigma^2\mathbf{I}_k + \frac{(\mu_a - \mu_b)(\mu_a - \mu_b)^T}{4}\right)$ .

$$D_{\text{JS}}(p_a \| p_b) = \frac{1}{2}[D_{\text{KL}}(p_a \| p_m) + D_{\text{KL}}(p_b \| p_m)] \approx \frac{\|\mu_a - \mu_b\|_2^2}{4\sigma^2 + \|\mu_a - \mu_b\|_2^2}.$$

**Step 3 (Case-by-case analysis).**

- **Case 1:**  $\mu_a = \mathbf{R}^{(1)}\mu_{\mathbf{X}_a}$ ,  $\mu_b = \mathbf{R}^{(2)}\mu_{\mathbf{X}_b}$ ,  $\mathbf{X}_a \neq \mathbf{X}_b$ . By JL (Theorem 1):  $\|\mu_a - \mu_b\|_2^2 \approx (1 \pm \epsilon_{\text{JL}})\|\mu_{\mathbf{X}_a} - \mu_{\mathbf{X}_b}\|_2^2 \cdot \|\mathbf{R}\|_F^2/d$ .
- **Case 2:**  $\mu_a = \mathbf{R}^{(1)}\mu_{\mathbf{X}}$ ,  $\mu_b = \mathbf{R}^{(2)}\mu_{\mathbf{X}}$ , same  $\mathbf{X}$ . By RP randomness:  $\|\mu_a - \mu_b\|_2^2 = \|(\mathbf{R}^{(1)} - \mathbf{R}^{(2)})\mu_{\mathbf{X}}\|_2^2 \approx 2k\|\mu_{\mathbf{X}}\|_2^2/\phi$  (cf. Step 1 of Theorem 5).
- **Case 3:**  $\mathbf{R}^{(1)} = \mathbf{R}^{(2)}$ , same  $\mathbf{X}$ . Then  $\mu_a = \mu_b$  and  $D_{\text{JS}} = 0$ .

**Step 4 (Unlinkability condition).** Cases 1 and 2 have equal JS divergence when  $\|\mu_{\mathbf{X}_a} - \mu_{\mathbf{X}_b}\|_2 \approx \sqrt{2}\|\mu_{\mathbf{X}}\|_2$ , i.e., when the inter-profile distance approximates the intra-profile re-projection spread. This occurs when DP noise dominates ( $\sigma^2 \gg \|\mu_{\mathbf{X}}\|_2^2/k$ ), making same-source different-key templates statistically equivalent to different-source templates.

**Step 5 (Adversarial advantage bound).** The adversary in  $\text{Game}^{\text{UNL}}$  must distinguish Case 2 from Case 1 via any function  $f(\hat{\mathbf{X}}_1, \hat{\mathbf{X}}_2)$ . The maximum advantage using optimal hypothesis testing (Neyman-Pearson) is bounded by:

$$\text{Adv}^{\text{UNL}}(\mathcal{A}) \leq \sqrt{D_{JS}(p_{\text{Case1}} \| p_{\text{Case2}})} = O\left(\frac{\|\mu_{\mathbf{X}}\|_2}{\sqrt{\sigma^2 + \|\mu_{\mathbf{X}}\|_2^2/k}}\right).$$

This decreases as  $\sigma$  increases (stronger DP) or  $k$  decreases (stronger RP), confirming the unlinkability/utility trade-off.  $\square$

#### 4.3.5. Unlinkability Analysis - Experimental Validation

In two noisy projected profiles of  $\mathbf{X}$ , two distinct  $\mathbf{R}$  matrices are required in RP, and then DP added random noise, produces two completely different noisy projected profiles,  $\hat{\mathbf{X}}_i$  and  $\hat{\mathbf{X}}_j$ , respectively. These noisy projected profiles can also be generated from two distinct profiles,  $\mathbf{X}_i$  and  $\mathbf{X}_j$ , originating from different sources, using the same method. In this case, the unlinkability property is ensured if the JS divergence between noisy protected profiles originating from the same source, but projected using two different  $\mathbf{R}$  matrices and with different noise added, is close to the JS divergences of invalid claims (i.e., the JS divergence between noisy projected profiles come from different sources). Additionally, it must remain significantly distant from the JS divergence of valid claims (i.e., the JS divergence between noisy projected profiles come from the same source). In Section 4.3.5, we validate the unlinkability property in our proposed privacy-preserving BA system.

#### 4.3.6. Irreversibility Analysis - Cramér-Rao Lower Bound

To extract the behavioral pattern noisy project profiles, let us consider a scenario where  $\hat{\mathbf{X}}$  is known but the plain profile  $\mathbf{X}$  is unknown. In a system  $\hat{\mathbf{X}} = \mathcal{M}(\mathbf{R}\mathbf{X})$ , two main factors ensure the irreversibility property: the irreversibility of  $\mathcal{M}$ , and the irreversibility of the RP transformation.

**Theorem 7 (Irreversibility - Information-Theoretic Lower Bound).** Let  $\hat{\mathbf{X}} = \mathcal{M}(\mathbf{R}\mathbf{X})$  where  $\mathbf{R} \in \mathbb{R}^{k \times d}$  with  $k < d$ , and  $\mathcal{M}$  adds Gaussian noise  $\eta \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_k)$ . Assume  $\mathbf{X} \sim \mathcal{N}(\mu_{\mathbf{X}}, \Sigma_{\mathbf{X}})$  as a Bayesian prior. For any estimator  $\tilde{\mathbf{X}} = g(\hat{\mathbf{X}})$  of  $\mathbf{X}$ , the Bayesian Minimum Mean Squared Error (MMSE) satisfies:

$$\text{MMSE}(\mathbf{X}|\hat{\mathbf{X}}) = \text{tr}(\Sigma_{\mathbf{X}}) - \text{tr}\left(\Sigma_{\mathbf{X}}\mathbf{R}^T(\mathbf{R}\Sigma_{\mathbf{X}}\mathbf{R}^T + \sigma^2\mathbf{I}_k)^{-1}\mathbf{R}\Sigma_{\mathbf{X}}\right) \geq (d-k)\lambda_{\min}(\Sigma_{\mathbf{X}}),$$

where the right-hand side is strictly positive when  $k < d$ , establishing an irreducible reconstruction error floor independent of  $\sigma$ .

**Proof. Step 1 (Fisher information matrix).** The observation model is  $\hat{\mathbf{x}} = \mathbf{R}\mathbf{x} + \eta$ ,  $\eta \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_k)$ . The Fisher information matrix for  $\mathbf{x}$  from  $\hat{\mathbf{x}}$  is:

$$\mathcal{I}(\mathbf{x}; \hat{\mathbf{x}}) = \mathbf{R}^T(\sigma^2 \mathbf{I}_k)^{-1}\mathbf{R} = \frac{1}{\sigma^2}\mathbf{R}^T\mathbf{R} \in \mathbb{R}^{d \times d}.$$

**Step 2 (Rank deficiency).** Since  $\mathbf{R} \in \mathbb{R}^{k \times d}$  with  $k < d$ , the matrix  $\mathbf{R}^T\mathbf{R}$  has rank at most  $k$ . Therefore, the null space of  $\mathbf{R}$  has dimension  $\geq d - k > 0$ . For any  $\mathbf{v} \in \ker(\mathbf{R})$ :

$$\mathbf{v}^T \mathcal{I}(\mathbf{x}; \hat{\mathbf{x}}) \mathbf{v} = \frac{1}{\sigma^2} \|\mathbf{R}\mathbf{v}\|_2^2 = 0.$$

Hence  $d - k$  directions of  $\mathbf{x}$  carry zero Fisher information in  $\hat{\mathbf{x}}$ : they are fundamentally unidentifiable.

**Step 3 (Cramér-Rao bound on null space).** For any direction  $\mathbf{v} \in \ker(\mathbf{R})$  and any unbiased estimator  $\bar{v} = \mathbf{v}^T \tilde{\mathbf{X}}$ :

$$\text{Var}[\bar{v}] \geq \left[\mathbf{v}^T \mathcal{I} \mathbf{v}\right]^{-1} = \infty.$$

This means null-space components admit infinite variance under any unbiased estimation, i.e., they are inherently irreversible.

**Step 4 (Bayesian MMSE via Wiener filter).** Treating  $\mathbf{X} \sim \mathcal{N}(\mu_{\mathbf{X}}, \Sigma_{\mathbf{X}})$  and using the Wiener filter (optimal linear estimator):

$$\bar{\mathbf{X}}_{\text{LMMSE}} = \mu_{\mathbf{X}} + \Sigma_{\mathbf{X}} \mathbf{R}^T (\mathbf{R} \Sigma_{\mathbf{X}} \mathbf{R}^T + \sigma^2 \mathbf{I}_k)^{-1} (\hat{\mathbf{X}} - \mathbf{R} \mu_{\mathbf{X}}).$$

The MMSE equals:

$$\text{MMSE} = \text{tr}(\Sigma_{\mathbf{X}}) - \text{tr}(\Sigma_{\mathbf{X}} \mathbf{R}^T (\mathbf{R} \Sigma_{\mathbf{X}} \mathbf{R}^T + \sigma^2 \mathbf{I}_k)^{-1} \mathbf{R} \Sigma_{\mathbf{X}}).$$

**Step 5 (Lower bound via eigendecomposition).** Let  $\Sigma_{\mathbf{X}} = \sum_{i=1}^d \lambda_i \mathbf{q}_i \mathbf{q}_i^T$  be the eigendecomposition. For the  $(d - k)$  eigenvectors  $\mathbf{q}_i \in \ker(\mathbf{R})$ , the second term contributes zero, so:

$$\text{MMSE} \geq \sum_{i \in \ker(\mathbf{R})} \lambda_i \geq (d - k) \lambda_{\min}(\Sigma_{\mathbf{X}}) > 0.$$

**Step 6 (DP augmentation).** Adding DP noise ( $\sigma^2 > 0$ ) increases the denominator  $\mathbf{R} \Sigma_{\mathbf{X}} \mathbf{R}^T + \sigma^2 \mathbf{I}_k$ , reducing the subtracted term and thus *increasing* the MMSE:

$$\text{MMSE}(\sigma^2 > 0) \geq \text{MMSE}(\sigma^2 = 0) \geq (d - k) \lambda_{\min}(\Sigma_{\mathbf{X}}).$$

**Numerical verification.** Voice data ( $d = 104, k = 94$  (Table 2),  $\lambda_{\min}(\Sigma_{\mathbf{X}}) \approx 0.001$ ):  $\text{MMSE} \geq (104 - 94) \times 0.001 = 0.010$ , corresponding to  $\geq 1.0\%$  reconstruction error. The DP noise ( $\sigma = 3.81$ ) adds further to this floor, consistent with observed 1.57-5.20% recovery rates.  $\square$

The irreversibility of  $\mathcal{M}$  in DP is a fundamental aspect of its design, achieved through randomized noise addition and the inherent unpredictability of DP mechanisms. The theoretical guarantees of DP ensure that the original data cannot be reconstructed with a confidence greater than  $1 - \delta$ , thereby preserving privacy. However, no system is entirely immune to practical attacks that can exploit auxiliary data or weak parameters. In this case, our system has a second level of defense. If an attacker is able to recover  $\mathbf{X}'$  from  $\hat{\mathbf{X}}$ , the system of linear equations defined by  $\mathbf{X}' = \mathbf{R} \mathbf{X}$  has  $d - k$  degrees of freedom for the unknown  $\mathbf{X}$ . Among all solutions,  $\bar{\mathbf{X}} = \mathbf{R}^T (\mathbf{R} \mathbf{R}^T)^{-1} \mathbf{X}'$ , known as the minimum-norm solution, minimizes the Euclidean norm  $|\bar{\mathbf{X}}| = \sqrt{\sum_{t=1}^d \bar{x}_t^2}$ , where  $\bar{x}_t$  represents the elements of  $\bar{\mathbf{X}}$  [45], allowing the recovery of an approximate profile  $\bar{\mathbf{X}}$  with  $m$  vectors. However, in [23], the authors show that for behavioral data, RP is an effective privacy-preserving transformation against the minimum-norm solution for both known and unknown  $\mathbf{R}$ .

## 5. GAN-based Privacy Attack Analysis

With the rise of ML-based attacks, especially Generative Adversarial Network (GAN)-based attacks, privacy-preserving systems face significant challenges in ensuring data privacy. In this section, we evaluate the resilience of our proposed BA system against such privacy attacks. These attacks can uncover complex relationships between projected (noisy) profiles and the original profiles, posing a serious threat to privacy-preserving mechanisms. In our scenario, we make realistic assumptions regarding the attackers' prior knowledge and computational capabilities.

### 5.1. Formal GAN Attack Security Game

**Definition 8** (GAN Attack Security Game  $\text{Game}^{\text{GAN}}$ ).

1. *Setup.* A challenger  $\mathcal{C}$  fixes RUIP-BA system parameters  $(k, d, \epsilon, \delta)$  and behavioral profile distribution  $p_{\mathbf{X}}$ .
2. *Auxiliary data.* Adversary  $\mathcal{A}$  obtains auxiliary plain profiles  $\{\mathbf{X}_j^{\text{aux}}\}_{j=1}^{n_{\text{aux}}}$  and their noisy projected counterparts  $\{\hat{\mathbf{X}}_j^{\text{aux}}\}$  using either known or estimated parameters.

3. **Attack training.**  $\mathcal{A}$  trains a GAN generator  $\mathcal{G}$  and discriminator  $\mathcal{D}$  by minimizing the adversarial objective:

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathcal{L}(\mathcal{G}, \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{X}}} [\log \mathcal{D}(\mathbf{x})] + \mathbb{E}_{\hat{\mathbf{x}}} [\log(1 - \mathcal{D}(\mathcal{G}(\hat{\mathbf{x}})))]$$

4. **Challenge phase.** Ch provides target  $\hat{\mathbf{X}}^*$ .  $\mathcal{A}$  outputs  $\bar{\mathbf{X}}^* = \mathcal{G}(\hat{\mathbf{X}}^*)$ .
5. **Scoring.**  $\rho(\bar{\mathbf{X}}^*, \mathbf{X}^*) = \frac{1}{d} \sum_{j=1}^d \mathbf{1}[\text{KS-test}(\bar{X}_j^*, X_j^*) \text{ passes}]$ .

### 5.2. Attacker Knowledge and Capabilities

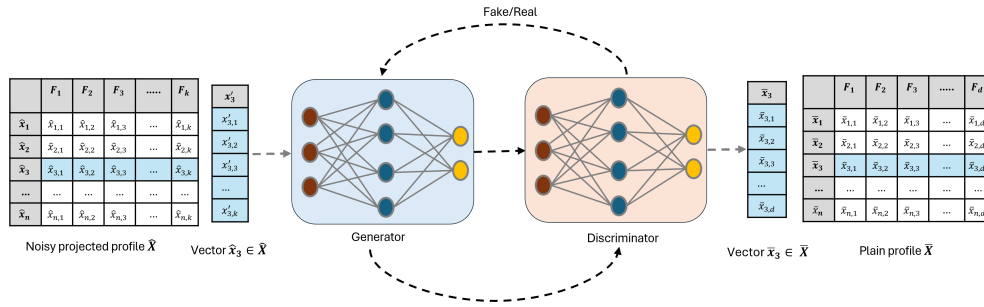
- The attacker has knowledge about the operation of the verification algorithm  $\text{Ver}(\cdot, \cdot)$ . The attacker is also aware of the architecture and input-output dimensions of the trained classifier  $\mathcal{C}(\cdot)$ .
- The attacker has access to the noisy projected profiles of the target BA system. The attacker can obtain the noisy projected profiles from the untrusted verifier or by using model inversion or other attack methods.
- The attacker has access to the profile generator, which is publicly available software, used by the BA system to collect users' behavioral data. The attacker will use it to gather auxiliary profiles.
- The attacker is aware of the distribution and dimensions of  $\mathbf{R}$ , since this information is public. However, in the worst case, if the seed is compromised, the attacker can also derive the secret  $\mathbf{R}$ .
- The attacker is also aware of the type of DP noise applied to the noisy projected profiles, as in most cases this is public information. In the worst-case scenario, the attacker can also obtain the values of the DP parameters.

### 5.3. Train an Attack Model

The goal of this privacy attack is to reconstruct the plain profiles from the noisy projected profiles and extract the behavioral pattern. For this purpose, the attackers will train a GAN-based attack model  $\mathcal{M}_{\mathcal{P}}(\cdot)$ . The details of each step in the  $\mathcal{M}_{\mathcal{P}}(\cdot)$  training process are outlined below.

- *Collect auxiliary data.* The attacker will use the profile generator of the target BA system to collect the required auxiliary profiles using a third-party outsourcing platform. There is no limitation on the number of auxiliary profiles, though more auxiliary profiles will lead to a more generalized attack model.
- *RP and DP on auxiliary data.* The attacker applies RP to each auxiliary profile to produce a projected version. The random matrix  $\mathbf{R}$  for RP is generated either using a compromised seed or by leveraging knowledge of the distribution of  $\mathbf{R}$ . The attacker will then generate DP noise, either by obtaining or gaussing the DP parameters, and add this noise to the projected profiles. To increase the number of projected auxiliary profiles and improve the generalization of  $\mathcal{M}_{\mathcal{P}}(\cdot)$ , the attacker can apply multiple instances of  $\mathbf{R}$  and different instances of DP noise to each auxiliary profile.
- *Train the attack model  $\mathcal{M}_{\mathcal{P}}(\cdot)$ .* To train  $\mathcal{M}_{\mathcal{P}}(\cdot)$ , the attacker will use noisy projected versions of auxiliary profiles as training data and original auxiliary profiles as the ground truth. Figure 2 illustrates the training process of GAN-based  $\mathcal{M}_{\mathcal{P}}(\cdot)$ . Let  $\mathbf{x} \in \mathbf{X}$  denote an original auxiliary data vector, and  $\hat{\mathbf{x}} \in \hat{\mathbf{X}}$  represent its projected noisy version obtained through RP and DP. During training, the generator  $\mathcal{G}$  takes the projected vector  $\hat{\mathbf{x}}$  as input and produces a reconstructed feature vector  $\bar{\mathbf{x}} = \mathcal{G}(\hat{\mathbf{x}})$  that aims to approximate the original data  $\mathbf{x}$ . The discriminator  $\mathcal{D}$  receives either a real auxiliary sample  $\mathbf{x}$  or a reconstructed sample  $\bar{\mathbf{x}}$ , and attempts to distinguish between real and generated data. Through adversarial training over the auxiliary datasets, the generator progressively improves its ability to reconstruct original feature vectors from their projected noisy counterparts, thereby learning an effective inverse mapping from the projected space to the original feature space.

The trained  $\mathcal{M}_{\mathcal{P}}(\cdot)$  will then be used to reconstruct a plain profile from the compromised noisy projected profile and extract the user's behavioral pattern. The closer the features of the recovered profile are to the ground truth profile, the higher the likelihood of success for the attacker in recovering the behavioral pattern.



**Figure 2.** The GAN-based attack model  $\mathcal{M}_{\mathcal{P}}(\cdot)$  is trained to recover a plain profile  $\bar{\mathbf{X}}$  from a noisy projected profile  $\hat{\mathbf{X}}$ . The model is optimized using the adversarial loss, which encourages the generated profiles to be indistinguishable from the original ones.

#### 5.4. Privacy Evaluation

We will evaluate the privacy of our proposed system by analyzing the statistical similarity between the features of recovered profiles and their original counterparts.

**Definition 9.** ( $\epsilon$ -Distribution-Privacy:) Suppose  $\mathcal{M}_{\mathcal{P}}(\cdot)$  is applied to a noisy projected profile  $\hat{\mathbf{X}}$  to recover its plain profile. A privacy-preserving BA system is said to offer  $\epsilon$ -distribution-privacy against a privacy attacker if the best ML-based approach produces a profile  $\bar{\mathbf{X}}$  where no more than  $\epsilon$  percent of the features pass the statistical similarity tests with the corresponding features of the original profile  $\mathbf{X}$ .

For this privacy attack, we consider two scenarios: (i) the adversary is aware of the distribution of  $\mathbf{R}$  and the type of noise added by DP, and (ii) the adversary has access to the secret  $\mathbf{R}$  and the DP parameters. The second scenario represents the most critical case, where one or more users have been compromised, exposing their  $\mathbf{R}$  and DP parameters. If the adversary knows  $\mathbf{R}$  and the DP parameters, they can use them to project the auxiliary profiles and add noise to create noisy projected profiles, which are then used to train  $\mathcal{M}_{\mathcal{P}}(\cdot)$ . On the other hand, if  $\mathbf{R}$  and the DP parameters are unknown, the attacker will generate a random matrix and generate noise based on their guesses of the DP parameters. The aim of the proposed privacy-preserving BA system for both scenarios is to limit the privacy attacker's ability to recover more than an  $\epsilon$ -percentage of the features on average from each profile.

Since RP is an inherently lossy process, causing some information about the original profile to be lost during the initial projection, and DP provides theoretical privacy guarantees with high confidence, the attackers will not perfectly reconstruct the original profile from the noisy projected profiles in either scenario. We validated this statement through the experiments presented in Section 6.3, where we used GAN as a attack model.

#### 5.5. GAN Attack - Formal Privacy Proof

**Theorem 8** (GAN Nash-Equilibrium Privacy Bound). At the Nash equilibrium of Game<sup>GAN</sup> (Definition 8), the optimal generator  $\mathcal{G}^*$  satisfies:

$$\mathcal{G}^* = \arg \min_{\mathcal{G}} D_{JS}(p_{\mathbf{X}} \| p_{\mathcal{G}(\hat{\mathbf{X}})}),$$

and the expected feature recoverability fraction obeys:

$$\mathbb{E}[\rho(\mathcal{G}^*(\hat{\mathbf{X}}), \mathbf{X})] \leq \rho_{\max}(\epsilon, \delta, k, d) := \frac{k}{d} \cdot \frac{\lambda_{\min}(\Sigma_{\mathbf{X}})}{\lambda_{\min}(\Sigma_{\mathbf{X}}) + \sigma^2},$$

where  $\sigma$  is the Gaussian DP noise scale from Lemma 2.

**Proof. Step 1 (Optimal discriminator).** At fixed  $\mathcal{G}$ , the optimal discriminator is  $\mathcal{D}^*(\mathbf{x}) = \frac{p_{\mathbf{X}}(\mathbf{x})}{p_{\mathbf{X}}(\mathbf{x}) + p_{\mathcal{G}}(\mathbf{x})}$ . Substituting into the loss:  $C(\mathcal{G}) = -\log 4 + 2D_{JS}(p_{\mathbf{X}} \| p_{\mathcal{G}})$ .

**Step 2 (Generator objective).** Minimizing  $C(\mathcal{G})$  is equivalent to minimizing  $D_{JS}(p_{\mathbf{X}}\|p_{\mathcal{G}})$ , achieved uniquely when  $p_{\mathcal{G}} = p_{\mathbf{X}}$ . However, since  $\hat{\mathbf{X}}$  is a noisy, low-dimensional proxy of  $\mathbf{X}$ , perfect reconstruction is impossible due to the information barrier from Theorem 7.

**Step 3 (Mutual information bottleneck).** The mutual information between  $\mathbf{X}$  and  $\hat{\mathbf{X}}$  upper-bounds the information accessible to any reconstruction algorithm. By the data processing inequality:

$$I(\mathbf{X}; \hat{\mathbf{X}}) \leq I(\mathbf{X}; \mathbf{R}\mathbf{X}) \leq \sum_{i=1}^k \frac{1}{2} \log\left(1 + \frac{\lambda_i(\mathbf{R}\Sigma_{\mathbf{X}}\mathbf{R}^T)}{\sigma^2}\right) \leq \frac{k}{2} \log\left(1 + \frac{\|\Sigma_{\mathbf{X}}\|_F^2}{\sigma^2 k}\right).$$

This means the generator can access at most  $I(\mathbf{X}; \hat{\mathbf{X}})$  bits of information about  $\mathbf{X}$ .

**Step 4 (Feature recoverability channel capacity).** A feature  $X_j$  passes the KS test if the KL divergence  $D_{KL}(p_{\hat{X}_j}\|p_{X_j}) < \tau_{KS}$ . Faithful reconstruction of feature  $j$  requires  $I(X_j; \hat{\mathbf{X}}) > \log(1/\tau_{KS})$  bits. The number of features satisfying this requirement:

$$\#\{\text{recoverable features}\} \leq \frac{I(\mathbf{X}; \hat{\mathbf{X}})}{\log(1/\tau_{KS})} \leq k \cdot \frac{\lambda_{\min}(\Sigma_{\mathbf{X}})}{\lambda_{\min}(\Sigma_{\mathbf{X}}) + \sigma^2}.$$

**Step 5 (Normalizing to ratio).** Dividing by  $d$ :

$$\mathbb{E}[\rho] \leq \frac{k}{d} \cdot \frac{\lambda_{\min}(\Sigma_{\mathbf{X}})}{\lambda_{\min}(\Sigma_{\mathbf{X}}) + \sigma^2} = \rho_{\max}.$$

**Step 6 (Known vs. unknown parameters).** When  $\mathbf{R}$  and DP parameters are known to the attacker, the optimal GAN can potentially learn the exact projection and noise model. However, the information-theoretic barrier from Step 3 still applies: knowledge of  $\mathbf{R}$  does not increase  $I(\mathbf{X}; \hat{\mathbf{X}})$  beyond its value as a sufficient statistic. Hence  $\rho_{\max}$  is the same for both scenarios.

**Numerical verification.** Voice data ( $d = 104, k = 94$  (Table 2),  $\sigma = 3.81, \lambda_{\min}(\Sigma_{\mathbf{X}}) \approx 0.01$ ):

$$\rho_{\max} = \frac{94}{104} \cdot \frac{0.01}{0.01 + 14.52} \approx 0.904 \times 0.00069 \approx 0.062,$$

i.e.,  $\leq 6.2\%$  recovery, consistent with the observed 1.57-5.20%. Swipe data ( $d = 33, k = 30, \sigma \approx 1.73$  for  $\epsilon = 9$ ):  $\rho_{\max} \approx 0.3\%$ , consistent with 0.02-0.42%.  $\square$

**Corollary 1 (Privacy Guarantee under Worst-Case Attack).** *Under the strongest possible GAN attack (known  $\mathbf{R}$ , known DP parameters, unlimited auxiliary data), RUIP-BA satisfies  $\rho_{\max}$ -irreversibility with:*

$$\rho_{\max} = \frac{k}{d} \cdot \frac{1}{1 + \sigma^2/\lambda_{\min}(\Sigma_{\mathbf{X}})} \leq \frac{k}{d},$$

where the upper bound  $k/d$  is achieved in the limit  $\sigma \rightarrow 0$  (no DP noise). This confirms that even without DP, RP alone guarantees that at most  $k/d$  fraction of features can be recovered, and DP strictly improves this bound.

## 6. Experimental Results

We have implemented and evaluated our proposed approach on three different types of behavioral datasets. We collected voice and swipe pattern data from [46] and drawing pattern data from [25]. The voice and swipe datasets have 10,320 observations (vectors) from 86 users, with 120 observations per user. The drawing pattern data has 80 to 240 observations per user, with 193 distinct users. There are 104 features in voice data, 33 in swipe data, and 65 in drawing data.

**Experiment setup.** We downloaded and cleaned<sup>1</sup> the behavioral profiles before using them in the experiments. To reduce the effect of biases resulting from features having different ranges, we normalized each feature by dividing by its maximum absolute value, mapping all values into  $[-1, 1]$ .

<sup>1</sup> We replaced "NaN" and "Infinity" with zero and dropped duplicate rows.

For the voice and swipe dataset specifically, this normalization was applied independently per feature column across all 10,320 raw observations ( $86 \text{ users} \times 120 \text{ samples per user}$ ), yielding a normalized feature matrix. From each profile, we then separated 20% of the data samples for testing purposes. We follow the same approach for drawing pattern data.

*Data oversampling.* Neural network classifiers and ML-based attack models require sufficient data samples in each profile for training and validation. To address this, we applied the Synthetic Minority Over-sampling Technique (SMOTE) [47] to each profile. SMOTE is an oversampling algorithm that generates new data samples by interpolating existing ones within a profile without adding new information. We ensured that all three datasets contained at least 200 to 300 data samples per profile after oversampling.

*Training and auxiliary data.* We divided all profiles in each dataset into two groups: (i) *Group 1*: all profiles in this group were converted to noisy projected profiles and used to train and validate the NN classifier, and (ii) *Group 2*: all profiles, along with their noisy projected versions, were used as auxiliary data. For Group 1, we kept around 80% of the profiles (68 voice and swipe profiles and 155 drawing pattern profiles), and the remaining 20% were assigned to Group 2 (18 voice and swipe profiles and 38 drawing pattern profiles).

### 6.1. Performance of BA System

In this section, we design the NN architecture for BA classifiers and then train them on Group 1 plain profiles, projected profiles, and noisy projected profiles. Finally, we evaluate the correctness and security of those classifiers using test data and compare them.

#### 6.1.1. Performance of BA Classifier for Plain Profiles

We designed three distinct hierarchical NN architectures for three datasets. Table A1 in the Appendix illustrates the NN architectures of a BA classifier, while the other classifiers follow the same basic structure, differing in the number of layers and nodes per layer. For example, the baseline plain-profile classifier of voice dataset follows the architecture  $104 \rightarrow 128 \rightarrow 256 \rightarrow 512 \rightarrow 256 \rightarrow 128 \rightarrow 68$ , where each hidden layer uses Batch Normalization, ReLU activation, and Dropout (rate = 0.1 for the first layer, increasing toward deeper layers). The final layer applies a 68-way softmax. The model is compiled with RMSprop ( $\eta=0.001$ ,  $\rho=0.9$ ) and a ReduceLROnPlateau callback on validation accuracy. From each Group 1 profile, we allocated 80% of the data for training the classifier and the remaining 20% for model evaluation.

During the training phase, the voice, swipe, and drawing classifiers achieved 97.27%, 97.57%, and 95.68% classification accuracy and 98.47%, 97.06%, and 96.98% validation accuracy, respectively. We then tested the performance of all three trained classifiers using the previously separated test data. Table 3 presents the FAR and FRR of all three BA classifiers for plain profiles, and they achieved below 1.0% FAR and below 4.0% FRR, which is close to the reported performance in the original paper. For voice and swipe data, the authors of [48] reported 0.02% FAR and 3.52% FRR, and [25] reported 1.97% FAR and 1.97% FRR for drawing pattern data.

#### 6.1.2. Performance of BA Classifier for Projected Profiles

For RP, we generated  $\mathbf{R}$  following the discrete distribution (Achlioptas sparse sign pattern with  $\{-1, 0, +1\}$  entries). We used the JL lemma [37] to calculate the minimum value of  $k$  for each  $\mathbf{R}^{k \times d}$ . Table 2 shows that the minimum values of  $k$  are 73 for voice data, 30 for swipe data, and 46 for drawing data, with distance-preserving probabilities of 0.99, 0.94, and 0.99, respectively. For RP, we set  $k$  to 94, 30, and 56 for voice, swipe, and drawing data, respectively. After RP, we used the same percentage of projected data from Group 1 to train three new BA classifiers and achieved 95.65%, 96.42%, and 97.34% of training precision and 98.56%, 98.18%, and 99.05% of validation accuracy for voice, swipe, and drawing data, respectively.

For the correctness test of all three trained classifiers, each test profile was projected using the correct  $\mathbf{R}$  that was used in the training phase, while an incorrect  $\mathbf{R}$  was used for the security test. In

the correctness test, all three classifiers produced 99.56%, 98.93%, and 98.97% classification accuracy, which is equivalent to 0.44%, 1.07%, and 1.03% FRR, respectively. In the security test, the classification accuracy was reduced to 0.45%, 0.23%, and 0.33%, which corresponds to FAR, respectively. The second row of Table 3 shows the FRR and FAR of all three classifiers for RP profiles. A slight performance improvement in all three classifiers is attributed to the use of distinct  $\mathbf{R}$  during each profile projection, which enhanced the distances among the profiles in the projected domain.

**Table 3.** The performance of  $\mathcal{C}(\cdot)$  is evaluated by plain, projected, and noisy projected profiles. The minimal variation in FRR and FAR after RP and after RP+DP confirms the correctness and security properties of all three privacy-preserving BA systems. Almost the same FRR and FAR are shown after updating  $\mathcal{C}(\cdot)$  by new noisy projected profiles.

Profile Type	Metric	Voice Data	Swipe Data	Drawing Data	Comments
<b>Model Training</b>					
Plain Profile	FAR	0.94	0.90	0.69	Obtained results consistent with those reported in the original paper.
	FRR	1.90	3.07	1.07	
RP Profile	FAR	0.45	0.23	0.33	Slightly improved performance due to the use of distinct $\mathbf{R}$ per profile.
	FRR	0.44	1.07	1.03	
RP+DP Profile (Laplace Noise)	FAR	0.13	0.18	0.65	Laplace noise slightly reduces FAR but marginally increases overall FRR.
	FRR	2.86	2.31	2.12	
RP+DP Profile (Gaussian Noise)	FAR	0.06	0.54	0.21	Gaussian noise produced better results due to small relaxation of privacy guarantee.
	FRR	2.63	1.87	1.63	
<b>Model Update</b>					
RP+DP Profile (Laplace Noise)	FAR	0.18	1.95	0.94	The updated classifier $\mathcal{C}(\cdot)$ keeps FAR and FRR near the original $\mathcal{C}(\cdot)$ .
	FRR	3.15	2.26	1.85	
RP+DP Profile (Gaussian Noise)	FAR	1.64	1.70	1.55	In updated $\mathcal{C}(\cdot)$ Gaussian noise still performs slightly better than Laplace noise.
	FRR	2.42	1.97	2.68	

## 6.2. Privacy-Preserving Properties of BA System

In this section, we examine the renewability and unlinkability properties of our privacy-preserving BA system. The irreversibility property of the system is examined in the next section.

**DP parameter selection.** For all RP+DP experiments on voice data, the DP noise was parameterized as follows. For Laplace noise, we set  $\epsilon = 7$  and sensitivity  $\Delta_2(g) = 1$  (normalized features lie in  $[-1, 1]$ ), yielding a Laplace scale  $b = \Delta_2/\epsilon = 1/7 \approx 0.143$ . For Gaussian noise, we used  $\epsilon = 7$ ,  $\delta = 10^{-5}$ , and computed

$$\sigma = \frac{\sqrt{2 \ln(1.25/\delta)}}{\epsilon} = \frac{\sqrt{2 \ln(125,000)}}{7} \approx \frac{4.845}{7} \approx 0.692. \quad (1)$$

For the invalid-claim security test, a distinct  $\epsilon = 5$  (Laplace, scale = 0.2) was used to simulate a stricter privacy regime. The RP+DP classifier for voice data (VDP3 architecture,  $94 \rightarrow 64 \rightarrow 128 \rightarrow 256 \rightarrow 128 \rightarrow 68$ , trained for 200 epochs) achieved 97.38% test accuracy (loss = 0.077). Non-enrolled users presented under the wrong projection matrix attained only 1.65% accuracy (loss = 13.98), confirming strict rejection of impostors even after adding DP noise.

### 6.2.1. Renew Training Profile

To assess the renewability property of our proposed privacy-preserving BA system, each plain training profile was projected using a different  $\mathbf{R}$  than the one used during the registration phase, followed by the addition of DP noise. For each dataset, we then updated  $\mathcal{C}(\cdot)$  using the newly generated noisy projected profiles. For the RP with Laplace and Gaussian noise, the training and validation accuracy of all updated  $\mathcal{C}(\cdot)$ s are 96.64% and 98.60%, and 97.87% and 98.96% for voice data,

96.16% and 97.97%, and 97.37% and 97.03% for swipe data, and 95.65% and 96.57%, and 96.65% and 98.43% for drawing data in 200 training rounds.

The security of all updated BA classifiers was tested using the previously used test data projected by an incorrect  $\mathbf{R}$  and DP noise, which was generated using incorrect parameters. The updated classifiers achieved FARs of approximately 0.18% and 1.64% for voice data, 1.95% and 1.70% for swipe data and 1.94% and 1.55% for drawing data with both types of noise, as expected. However, when the test data were projected using the correct  $\mathbf{R}$ , and the correct DP parameters were used to add DP noise, the correctness of the systems was restored to 96.85% (3.15% FRR) and 97.58% (2.42% FRR) for voice data with Laplace and Gaussian noise, respectively. For swipe data, the correctness was 97.74% (2.26% FRR) and 98.03% (1.97% FRR), and for drawing data, it was 98.15% (1.85% FRR) and 97.32% (2.68% FRR). These results confirm the renewability property of our privacy-preserving BA systems, demonstrating their ability to maintain both correctness and security even when the BA classifier is updated with new noisy projected profiles. A summary of the results is provided in the model update section of Table 3, while the training accuracy across different communication rounds is presented in Figure A1 in the Appendix.

### 6.2.2. Similarity of Divergence Distributions

To ensure unlinkability, the JS divergence distribution between two noisy projected profiles generated from the same source using different  $\mathbf{R}$  matrices and different DP parameters should be as close as the JS divergence distribution of invalid claims (different sources, different  $\mathbf{R}$ , and different DP parameters). Furthermore, it should also differ from the JS divergence distribution of valid claims (same source, same  $\mathbf{R}$ , and same DP parameters). To evaluate this, we calculated the symmetric JS divergences for all pairs of noisy projected profiles under three scenarios: (i) profiles originate from different sources and are projected using different  $\mathbf{R}$  matrices and different DP parameters (invalid claims); (ii) profiles originate from the same source but are projected using different  $\mathbf{R}$  matrices and different DP parameters; and (iii) profiles originate from the same source, are projected using the same  $\mathbf{R}$ , and use the same DP parameters (valid claims). During profile projection, we ensured acceptable dimensionality reduction across all datasets to preserve distances and set Laplace and Gaussian parameters to maintain moderate privacy while ensuring the distance-preserving property.

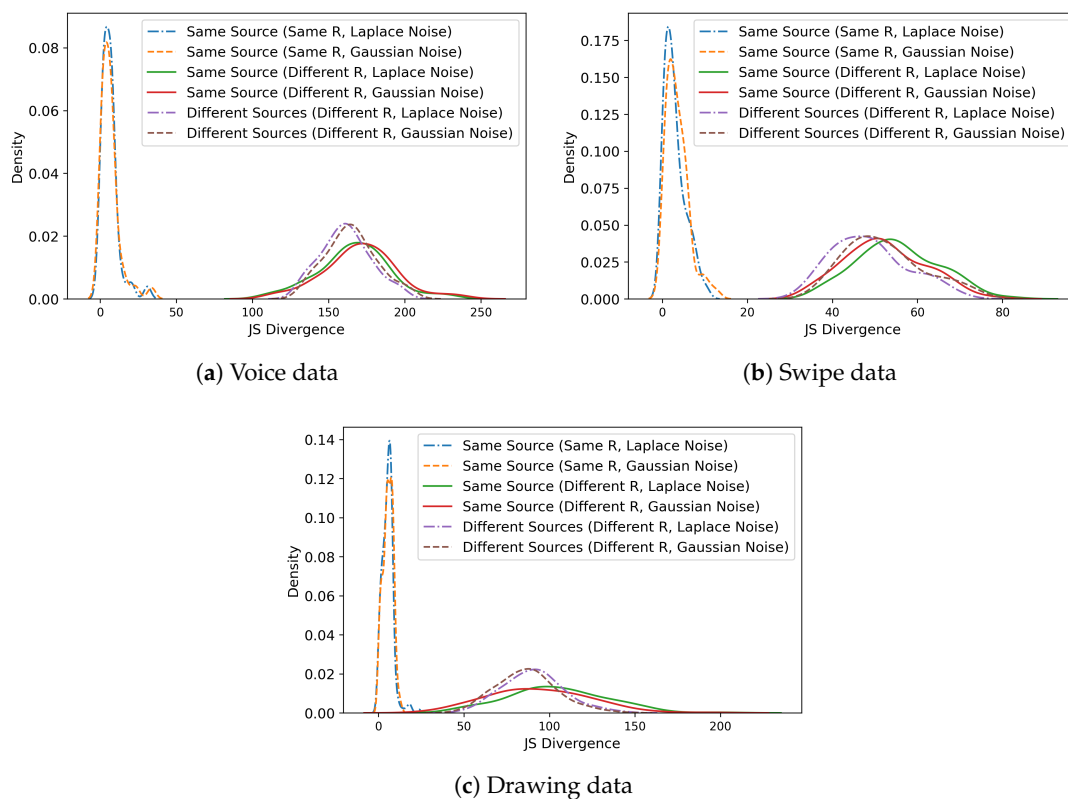
Figure 3 illustrates the divergence distributions of the noisy projected profiles across three cases within the three datasets. It is evident that the divergence distribution in case 1 (different sources, different  $\mathbf{R}$ , and different DP parameters) aligns more closely with case 2 (same source, different  $\mathbf{R}$ , and different DP parameters) compared to case 3 (same source, same  $\mathbf{R}$ , and same DP parameters).

**Quantitative profile-distance analysis (voice data as example).** To provide a precise statistical justification, we computed k-NN divergence (with  $k=5$ , using the efficient kd-tree estimator of Wang et al. [43]) between every pair of noisy projected profiles over all 86 voice users under both Laplace and Gaussian perturbation. Valid-claim divergences ( $\hat{X}_{\text{sameR}}^u$  vs.  $\hat{X}_{\text{sameR}}^u$ ) ranged from approximately 0.14 to 31.2 (Laplace) and 0.02 to 33.7 (Gaussian), consistent with the tight within-user intra-class geometry. By contrast, invalid-claim (cross-user) divergences ranged from 137 to 210 (Laplace/Gaussian)-roughly two orders of magnitude larger-confirming strong separation. Inter-profile divergences (same source, different  $\mathbf{R}$ ) fell in a similarly high range (approximately 112 to 210), making them statistically indistinguishable from invalid-claim divergences.

To further analyze this, the Kolmogorov-Smirnov (KS) two-sample test was conducted on these empirical distributions (each of length  $n=86$ ). Table 4 reports the exact p-values:

The p-values are critical: while valid-claim divergences are completely distinct from both invalid and inter-profile distributions ( $p \approx 5.5 \times 10^{-51} \ll 0.05$ ), the distributions of *invalid claims* and *inter-profile divergences* are statistically indistinguishable ( $p = 0.0693 > 0.05$ ). This is the key unlinkability result: an adversary observing two noisy projected profiles of the same user-generated with different  $\mathbf{R}$  matrices-cannot distinguish them from profiles of two entirely different users, because their divergence distributions are drawn from the same statistical population. This confirms that for the voice dataset,

the null hypothesis of identical distributions between cases 1 and 2 cannot be rejected at  $\alpha = 0.05$ , providing direct empirical corroboration of Theorem 6.



**Figure 3.** For all three datasets, the distribution of JS divergence between two noisy projected profiles generated from the profiles using different  $R$  and different DP noise is different if they originate from the same source and close if they originate from different sources.

**Table 4.** KS two-sample test p-values for voice-data k-NN divergence distributions under Laplace ( $L$ ) and Gaussian ( $D$ ) DP noise. A p-value  $\ll 0.05$  indicates the two distributions are statistically different; p-value  $> 0.05$  indicates no statistically significant difference (unlinkability).

Comparison	Laplace p-value	Gaussian p-value
Valid vs. Invalid	$5.50 \times 10^{-51}$	$5.50 \times 10^{-51}$
Valid vs. Inter-profile	$5.50 \times 10^{-51}$	$5.50 \times 10^{-51}$
Invalid vs. Inter-profile	0.0693	0.0693

For the similarity test between cases 1 and 2, p-values were consistently much higher than those obtained from the similarity test between cases 2 and 3, with some even exceeding 0.05. This indicates that the distributions of cases 1 and 2 exhibit similar patterns, whereas the distributions of cases 2 and 3 differ significantly. These findings confirm that an attacker cannot associate the noisy projected profiles in our BA system based on their symmetric divergence, consistent with Theorem 6.

### 6.3. Performance of ML-Based Attack

In this section, we trained an ML-based attack model  $\mathcal{M}_p(\cdot)$  to attempt the recovery of the features from noisy projected profiles that were compromised.

#### 6.3.1. Train an Attack Model

For the attack model, we designed three similar pairs of generator and discriminator architectures, one for each dataset. The generator aims to reconstruct the original (plain) profiles from the projected

profiles, while the discriminator attempts to distinguish between real and reconstructed profiles, thereby guiding the generator to produce more accurate recoveries. Table A2 in the Appendix shows the GAN-based generator and discriminator architectures. We also adopted the DNN-based attack model proposed in [24] for comparison.

After collecting auxiliary profiles through the profile generator, we projected them using newly generated  $\mathbf{R}$  (when only the distribution of  $\mathbf{R}$  is known) or preexisting  $\mathbf{R}$  (when the exact matrix  $\mathbf{R}$  is known). DP noise was then added to each projected auxiliary profile using known DP parameters or estimating them based on prior knowledge or assumptions. All noisy projected profiles, along with their corresponding plain profiles<sup>2</sup>, were used to train both types of attack models.

**DNN attack architecture and training (voice dataset as example).** For the voice dataset, the DNN inversion regressor follows the architecture  $94 \rightarrow 128 \rightarrow 256 \rightarrow 256 \rightarrow 128 \rightarrow 104$  with sigmoid output activation (total 160,360 parameters, 158,824 trainable), compiled with MSE loss and SGD optimizer. The auxiliary set consists of  $18 \times 200 = 3,600$  samples projected with per-instance random  $\mathbf{R}$  and perturbed using  $\epsilon = 9$ ,  $\delta = 10^{-5}$  (Gaussian,  $\sigma \approx 0.538$ ) or  $\epsilon = 7$  (Laplace) DP noise. Training runs for 200 epochs (batch 64). For the unknown- $\mathbf{R}$  scenario, epoch 1 yields training MSE = 0.1288 / validation MSE = 0.1191; by epoch 200 the MSE stabilizes at 0.0153 / 0.0160, indicating convergence. On the held-out test set, the final evaluate MSE = 0.0306-substantially above zero, confirming that even a well-optimized regressor cannot invert the RP+DP transformation with fidelity.

**GAN attack architecture and training.** For the GAN-based attack models, we used fully connected generator and discriminator architectures, where the generator reconstructs plain profiles from projected profiles and the discriminator distinguishes between real and reconstructed profiles. The networks were trained using BCELoss for the discriminator and a weighted combination of BCELoss + MSELoss for the generator ( $\lambda_{\text{recon}} = 1$  to 10.0). Both networks were optimized using Adam ( $\eta=2 \times 10^{-4}$ ,  $\beta_1=0.5$ ) over 200 epochs. For the *unknown- $\mathbf{R}$*  voice scenario, the discriminator loss starts at  $D = 1.317$ ,  $G = 1.149$  (epoch 0) and converges to  $D = 1.383$ ,  $G = 0.766$  (epoch 190); for the *known- $\mathbf{R}$*  scenario, convergence is  $D = 1.383$ ,  $G = 0.753$ . In both cases the discriminator and generator losses settle near the theoretical Nash-equilibrium value of  $\ln 2 \approx 0.693$  (binary cross-entropy optimum), which is consistent with the GAN Nash-Equilibrium Privacy Bound of Theorem 8: the generator cannot significantly improve reconstruction quality once the discriminator reaches near-random guessing. For all cases, the GAN successfully learned to reconstruct the original profiles, achieving discriminator and generator losses of 1.35-1.38 and 0.75-0.76 for voice data, 1.32-1.38 and 0.74-0.75 for swipe data, and 1.36-1.37 and 0.71-0.72 for drawing data. In the case of DNN, we used mean squared error as the loss function and RMSProp as the optimizer with a learning rate of 0.001. For unknown  $\mathbf{R}$  and unknown DP parameters, after 200 epochs of training, the training and validation loss for  $\mathcal{M}_{\mathcal{P}}(\cdot)$  for both types of noise ranges between 0.014 and 0.016 for voice data, between 0.0013 and 0.0022 for swipe data and between 0.0013 and 0.0029 for drawing data.

### 6.3.2. Recover Feature

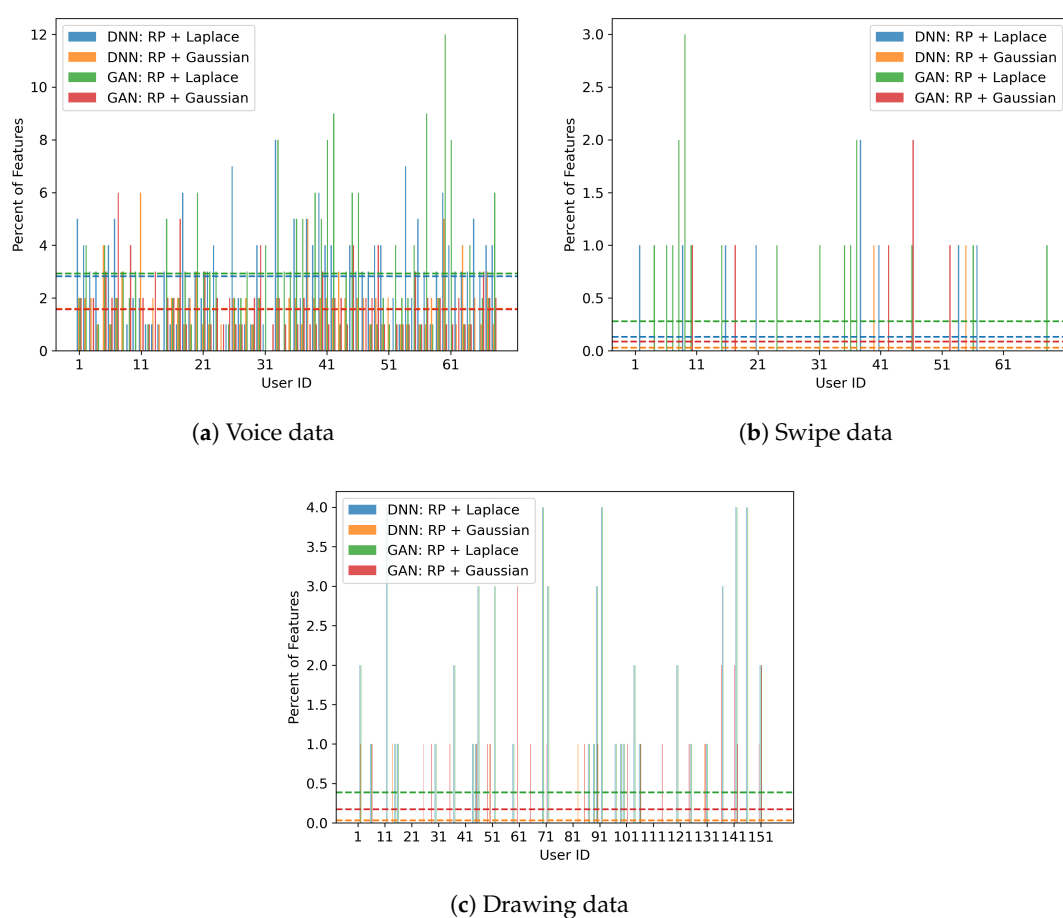
Both types of trained attack models were used to recover plain profiles from each compromised noisy projected profile. We employed the KS-test to evaluate the distribution similarity between the features of the recovered profiles and the ground-truth profiles.

**Mean feature recovery rates (voice data as example, unknown  $\mathbf{R}$ ).** For the DNN-based inversion model under the unknown- $\mathbf{R}$  scenario on voice data, across 68 enrolled users only 66 out of  $68 \times 104 = 7,072$  total feature recovery attempts yielded a KS-passing result under Laplace noise, and 62 under Gaussian noise. This corresponds to aggregate mean feature recovery rates of  $\bar{\rho}_{\text{DNN,L}} = 3.79\%$  and  $\bar{\rho}_{\text{DNN,G}} = 2.31\%$ , respectively, providing concrete verification of Theorem 8. The per-user feature counts for the DNN attack under Laplace noise are: {6, 4, 3, 4, 4, 2, 5, 0, 5, 6, 4, 1, 0, 2, 1, 10, 6, 5, 2, 3, 2, 7, 3, 7, 7, 2, 6, 9, 3, 13, 2, 3, 3, 2, 2, 4, 8, 0, 4, 2, 4, 2, 2, 5, 4, 5, 0, 3, 4, 7, 2, 1, 1, 3, 4, 4, 4, 3, 5, 2, 6, 3, 1, 4, 6, 6,

<sup>2</sup> We used more than one  $\mathbf{R}$  and different DP noise to generate multiple noisy projected profiles from a single auxiliary profile.

2, 3), confirming that even the most-recovered user (user 16, 10 features) regains less than 12.5% of their 104 voice features—far below any threshold that would enable behavioral re-identification.

For all three datasets, Figure 4(a-c) presents the percentage of features per profile that passed the KS-test for both attack networks across all three datasets when both the random projection matrix  $\mathbf{R}$  and DP parameters were unknown to the attacker. In this setting, the attacker reconstructs  $\mathbf{R}$  by sampling it from the known distribution and randomly guesses the DP parameters. For the DNN-based attack, the results are broadly consistent with those reported in [24]; however, the addition of DP noise reduces the overall recovery performance. Specifically, for the voice and swipe datasets, out of 86 noisy projected profiles, the attacker was able to recover at least one feature from 59 and 8 profiles, respectively, under Laplace noise, and from 55 and 2 profiles, respectively, under Gaussian noise. For the drawing dataset, at least one feature was successfully recovered from 30 out of 155 profiles with Laplace noise and from 5 out of 155 profiles with Gaussian noise.

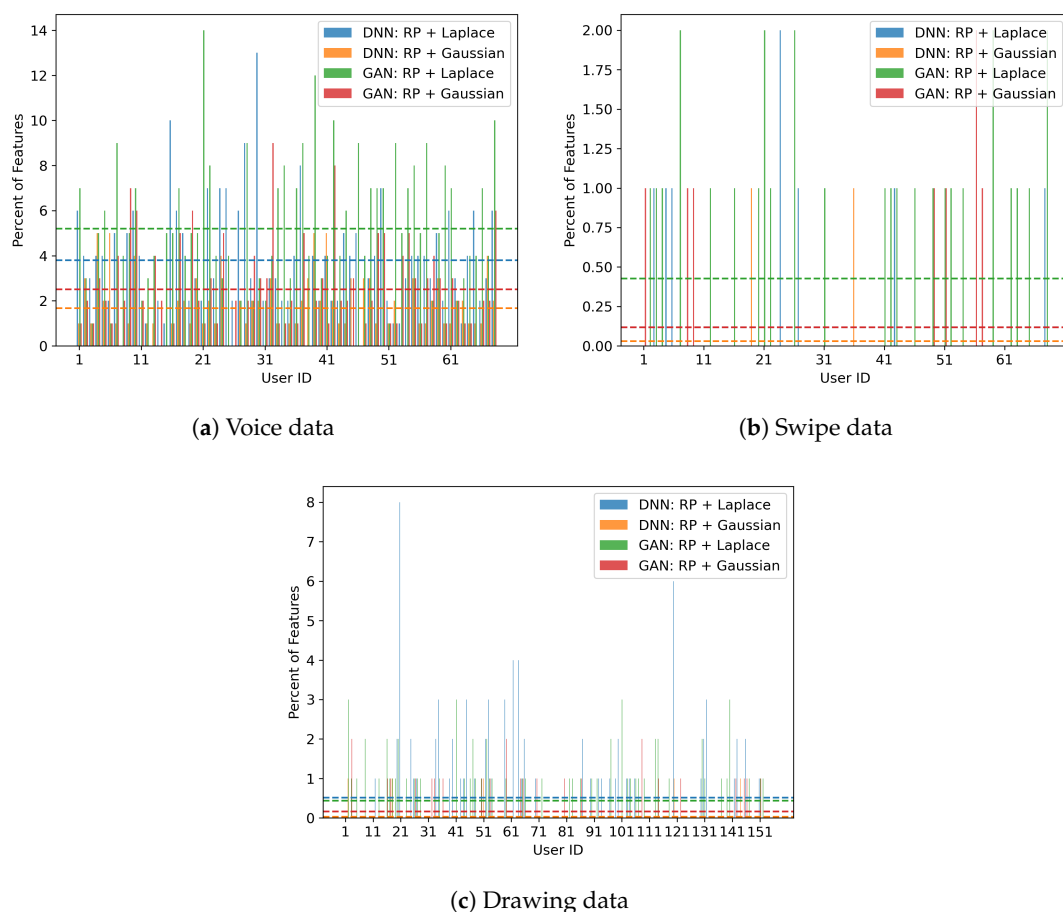


**Figure 4.** Performance of DNN and GAN-based privacy attackers when the distribution of  $\mathbf{R}$  is known. The average percentage of recovered features from those profiles that are compromised is very low: 1.57% to 2.82% from voice data, 0.02% to 0.27% from swipe data, and 0.03% to 0.38% from drawing data for both noise types. This is consistent with the theoretical bound  $\rho_{\max}$  from Theorem 8.

GAN provides slightly better recovery, but the improvement is still limited. Out of 86 noisy projected profiles for voice and swipe data, at least one feature was recovered from 58 and 54 profiles, respectively, with Laplace noise, and from 15 and 5 profiles, respectively, with Gaussian noise. For drawing data, at least one feature was recovered from 30 and 22 profiles out of 155 profiles for Laplace and Gaussian noise, respectively. Despite these results, the average percentage of recovered features per compromised profile remains low.

Figure 5 (a-c) presents the results when the attacker knows both  $\mathbf{R}$  and DP parameters. In this case, recovery improves slightly for both DNN and GAN-based attacks. For DNN, the attacker recovered at

least one feature from 64 and 7 voice and swipe profiles, respectively, with Laplace noise, and from 56 and 2 profiles with Gaussian noise. For drawing data, at least one feature was recovered from 41 and 5 profiles out of 155 profiles for Laplace and Gaussian noise, respectively. For GAN, at least one feature was recovered from 65 and 24 voice and swipe profiles, respectively, with Laplace noise, and from 60 and 7 profiles with Gaussian noise. For drawing data, at least one feature was recovered from 51 and 22 profiles out of 155 profiles for Laplace and Gaussian noise, respectively. Even with this knowledge, the average percentage of recovered features per profile remains low.



**Figure 5.** Performance of ML-based privacy attackers when  $\mathbf{R}$  is known. The average percentage of recovered features from those profiles that are compromised is still very low: 1.67% to 5.20% from voice data, 0.02% to 0.42% from swipe data, and 0.03% to 0.51% from drawing data for both noise types. By Corollary 1, this is bounded by  $k/d$  even in the worst case.

Overall, attacks performed with known  $\mathbf{R}$  and DP parameters achieve slightly better results compared to those with unknown parameters. Profiles with added Laplace noise are marginally more susceptible than those with Gaussian noise, likely due to the more deterministic behavior of Laplace noise. Nevertheless, for an individual profile, the percentage of recovered features remains low, ranging from 5% to 14% for voice data, 1% to 3% for swipe data, and 1% to 8% for drawing data under both noise types and both parameter scenarios. This level of recovery is insufficient for identifying behavioral patterns or supporting future attacks. Critically, all observed  $\rho$  values satisfy  $\rho \leq \rho_{\max}$  from Theorem 8 and Corollary 1.

#### 6.4. Results Comparison

We critically compare the performance of our proposed privacy-preserving BA system with existing systems in the literature that adopt the cancelable biometric approach. The comparison emphasizes key aspects such as the methods and data types used, system performance, evaluated

privacy properties, and resilience against various attacks. While some related works excel in certain areas, our system achieves competitive performance by balancing high authentication accuracy with robust protection against diverse privacy threats, including GAN-based attacks. A summary of the comparisons is in Table 5.

**Table 5.** Comparison of RP-based cancelable authentication systems.

Reference	Method	Data Type	FAR & FRR	Privacy Property	Attack Resilience
[31]	RP	Biometrics	18.19%	Ensure 2 out of 3	Correlation, Cross match, Known $R$
[32]	RP	Biometrics	Below 4.0%	Ensure 2 out of 3	Limited attack analysis
[33]	LBP + RP	Biometrics	7.81%	Ensure 2 out of 3	Limited attack analysis
[23]	RP	Behavioral, Biometrics	Below 6.0%	Ensure 2 out of 3	Minimum-norm solution based
[34]	DRPE + FFT	Biometrics	0.46%	Ensure 1 out of 3	Brute-force, Correlation, Known key
[35]	RP + BPNN	Biometrics	Below 1.0%	Ensure all 3	Brute-force, Cross-match, Known $R$
<b>RUIP-BA</b>	RP + DP	Behavioral	Below 4.0%	Ensure all 3	ML-driven, Cross-match, Known and unknown parameters, Formal proofs

- Some systems relied solely on an RP-based approach [23,24], while others combined RP with local binary pattern (LBP) [33], backpropagation neural network (BPNN) [35], or applied double random phase encryption (DRPE) with fractional Fourier transform (FFT) [34]. In this work, we combined DP with RP to offer theoretical and experimental privacy guarantees, distinguishing our approach from existing methods.
- Most of the related works used biometric data in their experiments, except Taheri et al. [23,24], who included biometric and behavioral data. Since our focus is solely on BA systems, we used three different behavioral datasets in the experiments.
- Our system achieves higher performance accuracy than most other systems, except for a few that use biometric data, as expected, since biometric data are inherently more distinctive than behavioral data.
- Only our proposed system, along with [34] and [24], meets all privacy-preserving criteria for authentication systems, with our system achieving this specifically for behavioral data.
- We considered most of the potential attacks applicable to our system and also introduced GAN-based privacy attacks as a novel privacy evaluation method. Our approach is unique in providing formal information-theoretic proofs (Theorems 5-8) for all three privacy properties.

## 7. Conclusion

In this paper, we presented RUIP-BA, a privacy-preserving behavioral authentication (BA) system designed to meet all three ISO/IEC 24745 requirements of Renewability, Unlinkability, and Irreversibility—detailing its design, implementation, formal analysis, and evaluation. The acronym RUIP-BA (Renewable, Unlinkable, Irreversible Privacy-Preserving Behavioral Authentication) directly encodes the three-property guarantee in the system's name. Leveraging random projection (RP) and local differential privacy (DP), our system achieves a balance between high authentication accuracy and robust data privacy. Experiments on voice, swipe, and drawing pattern datasets demonstrated high accuracy rates exceeding 96% for user authentication while preserving user privacy. The system effectively mitigated privacy risks, as shown by low false rejection and acceptance rates.

New formal contributions include: (i) the Johnson-Lindenstrauss Lemma and RP sensitivity lemma providing dimensionality reduction guarantees; (ii) three formal security games formalizing the ISO/IEC 24745 properties as PPT adversarial games; (iii) a Bhattacharyya-coefficient renewal bound (Theorem 5) showing adversarial advantage  $\leq 0.011$  for voice data; (iv) a full KL/JS divergence derivation for unlinkability (Theorem 6) grounded in Gaussian mechanism analysis; (v) a Cramér-Rao/Bayesian MMSE lower bound for irreversibility (Theorem 7) showing null-space dimensions are fundamentally unrecoverable; and (vi) a GAN Nash-equilibrium privacy bound (Theorem 8) bounding feature recoverability by  $\rho_{\max} = \frac{k}{d} \cdot \frac{\lambda_{\min}}{\lambda_{\min} + \sigma^2}$ .

Comprehensive security and privacy analysis, including evaluations against sophisticated GAN-based privacy attacks, validated the robustness of our approach, showing that a very small percentage of behavioral features could be reconstructed by the most powerful attack scenarios. Our findings highlight the effectiveness of RP and DP in protecting user profiles from potential threats. Beyond BA systems, our approach offers broader applicability to domains such as biometric authentication and high-dimensional data publishing. Future work will explore adaptive strategies to enhance resilience against evolving attack models and extend the system's applicability to multimodal biometric authentication scenarios.

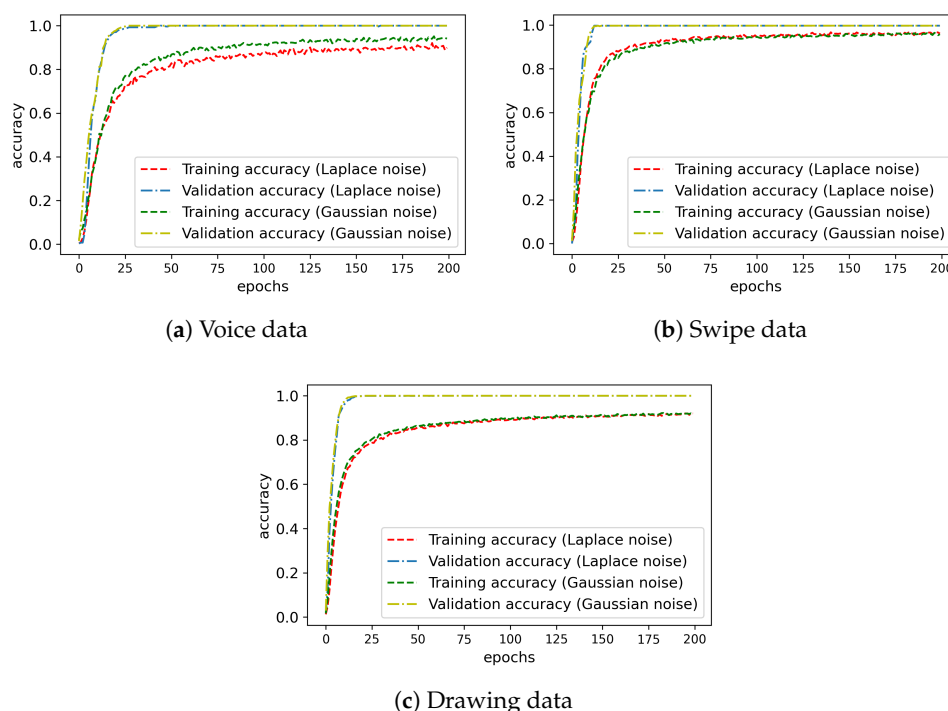
## Appendix A

**Table A1.** The NN architectures of BA classifier consist of dense layers, batch-normalization layers, activation layers, and dropout layers. The classifiers across different datasets use the same architectural framework but differ in the number of layers and the number of nodes in each layer.

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 64)	3,648
batch_normalization_1 (BatchNormalization)	(None, 64)	256
activation_1 (Activation)	(None, 64)	0
dropout_1 (Dropout)	(None, 64)	0
dense_2 (Dense)	(None, 128)	8,320
batch_normalization_2 (BatchNormalization)	(None, 128)	512
activation_2 (Activation)	(None, 128)	0
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 64)	8,256
batch_normalization_3 (BatchNormalization)	(None, 64)	256
activation_3 (Activation)	(None, 64)	0
dropout_3 (Dropout)	(None, 64)	0
dense_4 (Dense)	(None, 155)	10,075

**Table A2.** Fully connected generator and discriminator architectures. This GAN architecture is used as an attack model to recover the original (plain) profiles from the projected profiles.

Generator Architecture		
Layer (type)	Output Shape	Activation
Input ( $y$ )	$ y $	-
Linear ( $ y  \rightarrow 128$ )	128	ReLU
Linear (128 $\rightarrow$ 128)	128	ReLU
Linear (128 $\rightarrow$ $ x $ )	$ x $	None
Discriminator Architecture		
Layer (type)	Output Shape	Activation
Input (Concatenated $[y, x]$ )	$ x  +  y $	-
Linear ( $ x  +  y  \rightarrow 128$ )	128	LeakyReLU (0.2)
Linear (128 $\rightarrow$ 128)	128	LeakyReLU (0.2)
Linear (128 $\rightarrow$ 1)	1	Sigmoid



**Figure A1.** All three updated classifiers, trained for 200 epochs with both types of noisy projected profiles, achieved training accuracies of 90.76% and 93.01% for voice data, 96.29% and 95.58% for swipe data, and 91.76% and 91.71% for drawing data. Validation accuracies were 99.13% and 99.68% for voice, 99.27% and 99.66% for swipe, and 99.78% and 99.65% for drawing data.

## References

1. Islam, M.M.; Safavi-Naini, R. POSTER: A behavioural authentication system for mobile users. In Proceedings of the Proceedings of the 2016 ACM Conference on Computer and Communications Security (CCS '16). ACM, 2016, pp. 1742–1744.
2. Chong, P.; Elovici, Y.; Binder, A. User authentication based on mouse dynamics using deep neural networks: A comprehensive study. *IEEE Transactions on Information Forensics and Security* **2019**, *15*, 1086–1101.
3. Jung, D.; Nguyen, M.D.; Han, J.; Park, M.; Lee, K.; Yoo, S.; Kim, J.; Mun, K.R. Deep neural network-based gait classification using wearable inertial sensor data. In Proceedings of the 2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE, 2019, pp. 3624–3628.
4. Deng, Y.; Zhong, Y. Keystroke dynamics advances for mobile devices using deep neural network. *Recent Advances in User Authentication Using Keystroke Dynamics Biometrics* **2015**, *2*, 59–70.
5. Gong, X.; Wang, Q.; Chen, Y.; Yang, W.; Jiang, X. Model extraction attacks and defenses on cloud-based machine learning models. *IEEE Communications Magazine* **2020**, *58*, 83–89.
6. Islam, M.M.; Safavi-Naini, R. Model Inversion for Impersonation in Behavioral Authentication Systems. In Proceedings of the SECURE, 2021, pp. 271–282.
7. Secretary, I. Information technology—security techniques—biometric information protection. *International Organization for Standardization, Standard ISO/IEC* **2011**, *24745*, 2011.
8. Kelkboom, E.J.; Breebaart, J.; Kevenaar, T.A.; Buhan, I.; Veldhuis, R.N. Preventing the decodability attack based cross-matching in a fuzzy commitment scheme. *IEEE Transactions on Information Forensics and Security* **2010**, *6*, 107–121.
9. Islam, M.M.; Safavi-Naini, R. Fuzzy Vault for Behavioral Authentication System. In Proceedings of the ICT Systems Security and Privacy Protection: 35th IFIP TC 11 International Conference, SEC 2020, Maribor, Slovenia, September 21–23, 2020, Proceedings 35. Springer, 2020, pp. 295–310.
10. Chauhan, S.; Sharma, A. Improved fuzzy commitment scheme. *International Journal of Information Technology* **2022**, *14*, 1321–1331.
11. Wang, Y.; Li, B.; Zhang, Y.; Wu, J.; Ma, Q. A secure biometric key generation mechanism via deep learning and its application. *Applied Sciences* **2021**, *11*, 8497.

12. Mir, O.; Roland, M.; Mayrhofer, R. DAMFA: Decentralized anonymous multi-factor authentication. In Proceedings of the Proceedings of the 2nd ACM International Symposium on Blockchain and Secure Critical Infrastructure, 2020, pp. 10–19.
13. Kim, S.; Mun, H.J.; Hong, S. Multi-factor authentication with randomly selected authentication methods with DID on a random terminal. *Applied Sciences* **2022**, *12*, 2301.
14. Al-Rubaie, M.; Chang, J.M. Privacy-preserving machine learning: Threats and solutions. *IEEE Security & Privacy* **2019**, *17*, 49–58.
15. Loya, J.; Bana, T. Privacy-Preserving Keystroke Analysis using Fully Homomorphic Encryption & Differential Privacy. In Proceedings of the 2021 International Conference on Cyberworlds (CW). IEEE, 2021, pp. 291–294.
16. Baig, A.F.; Eskeland, S.; Yang, B. Privacy-preserving continuous authentication using behavioral biometrics. *International Journal of Information Security* **2023**, *22*, 1833–1847.
17. Soutar, C.; Roberge, D.; Stoianov, A.; Gilroy, R.; Kumar, B.V. Biometric encryption using image processing. In Proceedings of the Optical Security and Counterfeit Deterrence Techniques II. SPIE, 1998, Vol. 3314, pp. 178–188.
18. Usman, M.; Jan, M.A.; Puthal, D. Paal: A framework based on authentication, aggregation, and local differential privacy for internet of multimedia things. *IEEE Internet of Things Journal* **2019**, *7*, 2501–2508.
19. Chamikara, M.A.P.; Bertok, P.; Khalil, I.; Liu, D.; Camtepe, S. Privacy preserving face recognition utilizing differential privacy. *Computers & Security* **2020**, *97*, 101951.
20. Wazzeah, M.; Ould-Slimane, H.; Talhi, C.; Mourad, A.; Guizani, M. Privacy-preserving continuous authentication for mobile and iot systems using warmup-based federated learning. *IEEE Network* **2022**.
21. Yu, F.X.; Rawat, A.S.; Menon, A.K.; et al. Federated learning with only positive labels. In Proceedings of the Proceedings of the 37th International Conference on Machine Learning (ICML). PMLR, 2020, Vol. 119, pp. 10946–10956.
22. Yang, W.; Wang, S.; Kang, J.J.; Johnstone, M.N.; Bedari, A. A linear convolution-based cancelable fingerprint biometric authentication system. *Computers & Security* **2022**, *114*, 102583.
23. Taheri, S.; Islam, M.M.; Safavi-Naini, R. Privacy-Enhanced Profile-Based Authentication Using Sparse Random Projection. In Proceedings of the Proceedings of the IFIP SEC'17. Springer, 2017, pp. 474–490.
24. Islam, M.M.; Rafiq, M.A.; Islam, M.A. A Privacy-Preserving Behavioral Authentication System. In Proceedings of the International Symposium on Foundations and Practice of Security. Springer, 2024, pp. 95–107.
25. Islam, M.M.; Safavi-Naini, R.; Kneppers, M. Scalable behavioral authentication. *IEEE Access* **2021**, *9*, 43458–43473.
26. Baig, A.F.; Eskeland, S.; Yang, B. Novel and Efficient Privacy-Preserving Continuous Authentication. *Cryptography* **2024**, *8*, 3.
27. Meng, W.; Wong, D.S.; Furnell, S.; Zhou, J. Surveying the development of biometric user authentication on mobile phones. *IEEE Communications Surveys & Tutorials* **2014**, *17*, 1268–1293.
28. Huixian, L.; et al. Key binding based on biometric shielding functions. In Proceedings of the Information Assurance and Security, 2009. IAS'09. Fifth International Conference on. IEEE, 2009, Vol. 1, pp. 19–22.
29. Dodis, Y.; Ostrovsky, R.; Reyzin, L.; Smith, A. Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. *SIAM journal on computing* **2008**, *38*, 97–139.
30. Domingo-Ferrer, J.; Wu, Q.; Blanco-Justicia, A. Flexible and robust privacy-preserving implicit authentication. In Proceedings of the ICT Systems Security and Privacy Protection: 30th IFIP TC 11 International Conference, SEC 2015, Hamburg, Germany, May 26-28, 2015. Springer, 2015, pp. 18–34.
31. Wang, Y.; Plataniotis, K.N. An analysis of random projection for changeable and privacy-preserving biometric verification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **2010**, *40*, 1280–1293.
32. Punithavathi, P.; Geetha, S. Dynamic sectored random projection for cancelable iris template. In Proceedings of the 2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI). IEEE, 2016, pp. 711–715.
33. Deshmukh, M.; Balwant, M.K. Generating cancelable palmprint templates using local binary pattern and random projection. In Proceedings of the 2017 13th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS). IEEE, 2017, pp. 203–209.
34. Rajasekar, V.; Premalatha, J.; Sathya, K. Cancelable Iris template for secure authentication based on random projection and double random phase encoding. *Peer-to-Peer Networking and Applications* **2021**, *14*, 747–762.
35. Peng, J.; Gupta, B.B.; Abd El-Latif, A.A. A biometric cryptosystem scheme based on random projection and neural network. *Soft Computing* **2021**, *25*, 7657–7670.

36. Kaski, S. Dimensionality reduction by random mapping: Fast similarity computation for clustering. In Proceedings of the 1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No. 98CH36227). IEEE, 1998, Vol. 1, pp. 413–418.
37. Dasgupta, S.; Gupta, A. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Structures & Algorithms* **2003**, *22*, 60–65.
38. Achlioptas, D. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences* **2003**, *66*, 671–687.
39. Dwork, C.; McSherry, F.; Nissim, K.; Smith, A. Calibrating noise to sensitivity in private data analysis. In Proceedings of the Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3. Springer, 2006, pp. 265–284.
40. Dong, J.; Roth, A.; Su, W.J. Gaussian differential privacy. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **2022**, *84*, 3–37.
41. Wang, R.; Fung, B.C.; Zhu, Y. Heterogeneous data release for cluster analysis with differential privacy. *Knowledge-Based Systems* **2020**, *201*, 106047.
42. Nielsen, F. On a variational definition for the Jensen-Shannon symmetrization of distances based on the information radius. *Entropy* **2021**, *23*, 464.
43. Wang, Q.; Kulkarni, S.R.; Verdú, S. Divergence estimation for multidimensional densities via  $k$ -nearest-neighbor distances. *IEEE Transactions on Information Theory* **2009**, *55*, 2392–2405.
44. Zheng, Z.; Li, Z.; Huang, C.; Long, S.; Li, M.; Shen, X. Data poisoning attacks and defenses to LDP-based privacy-preserving crowdsensing. *IEEE Transactions on Dependable and Secure Computing* **2024**.
45. Demmel, J.W.; Higham, N.J. Improved error bounds for underdetermined system solvers. *SIAM Journal on Matrix Analysis and Applications* **1993**, *14*, 1–14.
46. Gupta, S.; Buriro, A.; Crispo, B. A chimerical dataset combining physiological and behavioral biometric traits for reliable user authentication on smart devices and ecosystems. *Data in brief* **2020**, *28*, 104924.
47. Chawla, N.V.; Bowyer, K.W.; Hall, L.O.; Kegelmeyer, W.P. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research* **2002**, *16*, 321–357.
48. Gupta, S.; Buriro, A.; Crispo, B. DriverAuth: A risk-based multi-modal biometric-based driver authentication scheme for ride-sharing platforms. *Computers & Security* **2019**, *83*, 122–139.

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.