



HAL
open science

Privacy-preserving Biometric Authentication Systems, a Cryptographic Approach

Axel Durbet

► **To cite this version:**

Axel Durbet. Privacy-preserving Biometric Authentication Systems, a Cryptographic Approach. Cryptography and Security [cs.CR]. Université Clermont Auvergne, 2024. English. NNT : 2024UCFA0113 . tel-04906100

HAL Id: tel-04906100

<https://theses.hal.science/tel-04906100v1>

Submitted on 22 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

UNIVERSITÉ CLERMONT AUVERGNE

École doctorale **Science Pour L'Ingénieur (SPI)**

Unité de recherche **Laboratoire d'Informatique, de Modélisation et d'Optimisation des Systèmes**

Thèse présentée par **Axel DURBET**

Soutenue le **5 novembre 2024**

En vue de l'obtention du grade de docteur de l'Université Clermont Auvergne

Discipline **Informatique**

Une approche cryptographique des systèmes d'authentification biométrique respectant la vie privée

Composition du jury

| | | | |
|----------------------------|------------------------|--|-------------------|
| <i>Rapporteurs</i> | Patrick LACHARME | Maître de Conférence HDR à l'Université de Caen Normandie | |
| | Geoffroy COUTEAU | Chargé de recherche HDR au CNRS | |
| <i>Examineurs</i> | Alexis BONNECAZE | Professeur à l'Aix-Marseille Université (AMU) | président du jury |
| | Reihaneh SAFAVI-NAINI | Professeure à l'University of Calgary | |
| | Sonia BEN MOKHTAR | Directrice de recherche au CNRS | |
| <i>Invité</i> | Koray KARABINA | Directeur de recherche au CNRC | |
| <i>Directeurs de thèse</i> | Pascal LAFOURCADE | Professeur à l'Université Clermont Auvergne (UCA) | |
| | Kévin THIRY-ATIGHEHCHI | Maître de Conférence à l'Université Clermont Auvergne(UCA) | |
| | Paul-Marie GROLLEMUND | Maître de Conférence à l'Université Clermont Auvergne(UCA) | |

UNIVERSITÉ CLERMONT AUVERGNE

Doctoral School **Science Pour L'Ingénieur (SPI)**

University Department **Laboratoire d'Informatique, de Modélisation et d'Optimisation des Systèmes**

Thesis defended by **Axel DURBET**

Defended on **5th November, 2024**

In order to become Doctor from Université Clermont Auvergne

Academic Field **Computer Science**

Privacy-preserving Biometric Authentication Systems, a Cryptographic Approach

Committee members

| | | | |
|--------------------|------------------------|--|---------------------|
| <i>Referees</i> | Patrick LACHARME | HDR Associate Professor at Université de Caen Normandie | |
| | Geoffroy COUTEAU | HDR Junior Researcher at CNRS | |
| <i>Examiners</i> | Alexis BONNECAZE | Professor at Aix-Marseille Université (AMU) | Committee President |
| | Reihaneh SAFAVI-NAINI | Professor at University of Calgary | |
| | Sonia BEN MOKHTAR | Senior Researcher at CNRS | |
| <i>Guest</i> | Koray KARABINA | Senior Researcher at CNRC | |
| <i>Supervisors</i> | Pascal LAFOURCADE | Professor at Université Clermont Auvergne (UCA) | |
| | Kévin THIRY-ATIGHEHCHI | Associate Professor at Université Clermont Auvergne(UCA) | |
| | Paul-Marie GROLLEMUND | Associate Professor at Université Clermont Auvergne(UCA) | |

UNE APPROCHE CRYPTOGRAPHIQUE DES SYSTÈMES D'AUTHENTIFICATION BIOMÉTRIQUE RESPECTANT LA VIE PRIVÉE

Résumé

Cette thèse vise à mettre en évidence les vulnérabilités des systèmes biométriques et à proposer des solutions pour renforcer la sécurité de ces données.

Un système biométrique permet d'authentifier ou d'identifier un individu en utilisant des caractéristiques physiques ou comportementales, telles que les empreintes digitales. Pour des raisons de sécurité, ces données ne sont pas utilisées en clair mais sont transformées en gabarits, rendant difficile la reconstitution des données originales. Cette transformation assure le respect de la vie privée des individus tout en permettant une authentification et une identification précises. En raison de leur utilisation analogue, nous avons étudié les données biométriques de manière similaire aux mots de passe en cryptographie. Plus précisément, nous avons d'abord étudié la probabilité d'occurrence d'une quasi-collision, c'est-à-dire la probabilité que deux gabarits de deux utilisateurs distincts soient proches. Les quasi-collisions posent problème car elles dégradent la capacité de reconnaissance du système et peuvent être exploitées par un attaquant cherchant à usurper l'identité de plusieurs utilisateurs. Pour éviter ces inconvénients, nous avons établi une borne sur la taille de la base de données pour prévenir les quasi-collisions et introduit un score pour aider à paramétrer les algorithmes de reconnaissance biométrique.

Ensuite, nous avons étudié les attaques par recherche exhaustive sur les données biométriques. Nous avons d'abord examiné les attaques ciblées, visant un utilisateur particulier, en étudiant la probabilité qu'un attaquant réussisse à usurper l'identité d'un utilisateur choisi dans différents scénarios. Cette étude nous a permis de définir des bornes de sécurité pour les bases de données de gabarits et de fournir des recommandations concernant les paramètres de sécurité pour les systèmes biométriques. Nous avons également investigué les attaques non-ciblées, où l'attaquant ne vise aucun utilisateur en particulier, pour évaluer la probabilité qu'un ou plusieurs attaquants réussissent à usurper l'identité de quelqu'un dans une base de données. Même si la probabilité d'usurper l'identité d'un individu spécifique est faible, il peut être facile d'usurper l'identité de quelqu'un lorsque la base de données est grande, de la même manière qu'il est probable que "0000" soit le mot de passe de quelqu'un dans une grande base de données. Cette analyse nous a permis de compléter notre investigation de la sécurité des données biométriques et de caractériser leurs limites.

Les attaques mentionnées ci-dessus s'appliquent principalement hors ligne. En ligne, ces attaques sont généralement détectées ou des contre-mesures sont mises en place pour ralentir les attaquants, comme l'augmentation du temps d'attente entre les tentatives. Pour limiter les problèmes liés aux attaques hors ligne, nous avons développé deux nouveaux protocoles d'authentification biométrique résistants aux attaques hors ligne. Le premier protocole utilise une preuve à divulgation nulle de connaissances pour garantir qu'un client malveillant, même avec des ressources de calcul illimitées, ne puisse obtenir aucune information utile du serveur pour effectuer une recherche exhaustive hors ligne. Le second protocole permet de corriger la donnée biométrique fournie par le client, conçu de telle sorte qu'un client malveillant avec une capacité de calcul polynomiale ne puisse obtenir aucune information utile.

Mots clés : sécurité biométrique ; transformations biométriques ; authentification biométrique ; identification biométrique ; extracteur flou réutilisable ; esquisse sécurisée réutilisable ; extracteur flou opaque ; esquisse sécurisée opaque ; calcul sécurisé ; preuve par jeu ; divulgation nulle de connaissances ; distance obfusquée ; correspondance floue ; distance de hamming ; fuite d'information ; problème du collecteur de coupons ; problème de la chaîne la plus proche ; quasi-collisions

PRIVACY-PRESERVING BIOMETRIC AUTHENTICATION SYSTEMS, A CRYPTOGRAPHIC APPROACH**Abstract**

This thesis aims to identify vulnerabilities in biometric systems and propose solutions to enhance their security. A biometric system authenticates or identifies an individual using physical or behavioral characteristics, such as fingerprints. For security reasons, these data are transformed into templates, making it difficult to reconstruct the original data. This transformation ensures privacy while enabling accurate authentication and identification. Due to their similar usage, we studied biometric data similarly to cryptographic passwords.

First, we examined the probability of near-collisions, where two templates from different users are similar. Near-collisions are problematic as they degrade system recognition and can be exploited by attackers to impersonate multiple users. To mitigate this, we established a database size limit to prevent near-collisions and introduced a score to help configure biometric recognition algorithms.

Next, we studied exhaustive search attacks on biometric data. We first focused on targeted attacks, which aim at a specific user, analyzing the probability of an attacker successfully impersonating a chosen user under various scenarios. This allowed us to define security bounds for template databases and provide recommendations for biometric system security parameters. We also investigated untargeted attacks, where the attacker does not aim at any specific user, to evaluate the probability of one or more attackers successfully impersonating someone in a database. Even if the probability of impersonating a specific individual is low, it can be easier to impersonate someone in a large database, similar to how "0000" is likely to be someone's PIN in a large dataset. This analysis completed our investigation of biometric data security and the characterization of their limitations.

The attacks described above are primarily offline. Online, these attacks are generally detected or countermeasures are implemented to slow down attackers, such as increasing the waiting time between two attempts. To address offline attack issues, we developed two new biometric authentication protocols resistant to offline attacks. The first protocol uses a zero-knowledge proof to ensure that a malicious client, even with unlimited computational resources, cannot obtain useful information from the server to perform an offline exhaustive search. The second protocol allows for the correction of biometric data provided by the client, designed so that a malicious client with polynomial computational capacity cannot obtain useful information.

Keywords: biometric security; biometric transformations; biometric authentication; biometric identification; reusable secure sketch; oblivious sketch; reusable fuzzy extractor; oblivious fuzzy extractor; secure computation; game-based; zero-knowledge; obfuscated distance; fuzzy matcher; hamming distance; information leakage; coupon collector problem; closest-string problem; near-collisions

Remerciements

Je tiens à exprimer ma profonde gratitude à toutes les personnes qui ont contribué à la réalisation de cette thèse.

En premier lieu, je remercie chaleureusement mon directeur de thèse, Pascal Lafourcade, Professeur au LIMOS et mes encadrants Kévin Thiry-Atighehchi, Maître de conférences au LIMOS, et Paul-Marie Grollemund, Maître de conférences au LMBP, pour leur encadrement exceptionnel, leurs conseils avisés, et leur soutien indéfectible tout au long de ce travail. Leur rigueur scientifique et leur enthousiasme ont été une source constante de motivation et d'inspiration.

Je tiens également à remercier les rapporteurs de cette thèse, Patrick Lacharme, HDR et Maître de conférences au GREYC, et Geoffroy Couteau, HDR et Chargé de recherche au CNRS, pour la relecture de cette thèse. Je suis également reconnaissant envers les examinateurs, Alexis Bonneau, Professeur à l'I2M, Reihaneh Safavi-Naini, Professeur à l'Université de Calgary, et Sonia Ben Mokhtar, Directrice de recherche au CNRS, pour avoir accepté de juger ce travail.

Un grand merci à Koray Karabina Directeur de Recherche au CNRS pour nos discussions stimulantes, ses encouragements, et son amitié. Travailler avec toi a été une expérience enrichissante tant sur le plan professionnel que personnel.

Je remercie chaleureusement mes collègues du LIMOS et de l'IUT d'Aurillac pour l'environnement de travail agréable et collaboratif. Vos encouragements et votre soutien m'ont beaucoup aidé.

Un merci particulier à mes collègues doctorants, Dorine Chagnon et Togo Jean Yves Kioye, pour leur soutien, leur camaraderie et leurs discussions enrichissantes. Travailler à vos côtés a été une expérience extrêmement précieuse.

Je tiens aussi à remercier les stagiaires qui ont contribué à nos travaux et à la vie du laboratoire sur le site délocalisé.

Je tiens aussi à exprimer ma reconnaissance à mes collègues enseignants, Clément Jacq, Denis Migdal et Maeva Paradis, pour leurs conseils et les échanges fructueux.

Un grand merci à mes amis et ma famille pour leur soutien inconditionnel et leur patience tout au long de ces années. Votre compréhension et vos encouragements m'ont permis de surmonter les moments difficiles.

À ma compagne, je souhaite exprimer ma plus profonde gratitude pour sa patience, et son soutien inébranlable tout au long de cette aventure. Merci d'avoir été à mes côtés dans les moments de doute et de célébrer chaque petite victoire avec moi.

Je remercie également l'Agence Nationale de la Recherche (ANR) pour son soutien, dans le cadre de la subvention ANR-20-CE39-0005 (projet PRIVABIO) pour leur soutien financier, sans lequel ce projet n'aurait pas été possible.

À toutes les personnes qui, de près ou de loin, ont contribué à la réalisation de cette thèse, je vous adresse mes remerciements les plus sincères.

Contents

| | |
|--|-----------|
| 1 Introduction | 1 |
| 1.1 About Biometrics | 2 |
| 1.1.1 History of Biometric | 2 |
| 1.1.2 Biometric Definition | 4 |
| 1.2 About the Data | 4 |
| 1.2.1 Personal Data | 4 |
| 1.2.2 Biometric Data | 5 |
| 1.2.3 Comparison Between Personal Data and Biometric Data | 5 |
| 1.3 Biometric Authentication Systems | 5 |
| 1.3.1 Specification | 6 |
| 1.3.2 Accuracy Metrics | 7 |
| 1.3.3 The Threat Points of a Biometric System | 9 |
| 1.3.4 Typology of Biometric Attacks | 10 |
| 1.4 Cryptographic Methods for Fuzzy Data | 11 |
| 1.5 Biometric System Problems Addressed | 13 |
| 1.5.1 Exhaustive Search Attacks Complexity | 14 |
| 1.5.2 Design of a Remote Secure Sketch Resistant to Offline Attacks | 14 |
| 1.6 Publications, and Submitted Works | 14 |
| 1.6.1 Attacks on Biometric Systems | 14 |
| 1.6.2 Construction of Cryptographic Primitive | 15 |
| 1.6.3 Related Publications | 15 |
| 2 Security of Biometrics Systems | 17 |
| 2.1 Introduction | 18 |
| 2.2 Master Templates and Near-Collisions | 21 |
| 2.2.1 Definitions | 21 |
| 2.2.2 Database Partitioning | 22 |
| 2.2.3 Multi-Near-Collisions and Master-templates | 31 |
| 2.2.4 Risk related to Master-templates | 36 |
| 2.2.5 Numerical Evaluations: Databases Security w.r.t Near-Collisions | 38 |
| 2.2.6 Metric Based Analysis | 38 |
| 2.3 Targeted and Untargeted Attack Models | 42 |
| 2.4 Targeted Attack | 43 |
| 2.4.1 Exhaustive Search | 43 |
| 2.4.2 The Center Search Attack | 45 |
| 2.4.3 The Impact of Information Leakage on Targeted Attack | 46 |
| 2.4.4 Exploiting the Leakage | 48 |
| 2.4.5 Accumulation Attack: A Passive Attack | 50 |
| 2.4.6 Weaker Attack Model: Compromise a point in the ball | 52 |
| 2.5 Untargeted Attack | 52 |
| 2.5.1 Metric Space Based Bounds | 52 |
| 2.5.2 Biometric security: Assess the Security Score Against the Outsider | 59 |
| 2.6 Conclusion and Future Work | 60 |

| | |
|--|------------|
| 3 Remote Secure Sketch | 63 |
| 3.1 Introduction | 64 |
| 3.2 Generic Construction of Secure Sketches from Groups | 67 |
| 3.2.1 Mathematical Construction: Groups with Unique t -factorization | 67 |
| 3.2.2 Computational Secure Sketch Construction | 70 |
| 3.2.3 Security of SPS | 72 |
| 3.2.4 On the Hardness of and Relationships between DSPP and CSPP | 75 |
| 3.2.5 Robust SPS in the Random Oracle Model | 84 |
| 3.2.6 Concrete Instantiations of SPS | 85 |
| 3.2.7 Details on VBB and NTT | 85 |
| 3.2.8 Robust Fuzzy Extractor based on Robust SPS | 88 |
| 3.3 Biometric Authentication Protocol from Groups | 89 |
| 3.4 Remote Secure Sketch from Group | 92 |
| 3.4.1 Requirements for Building a RSS from Subset Product Problems | 93 |
| 3.4.2 Mechanics of our Remote Subset Product Sketch | 93 |
| 3.4.3 Authentication and Secrecy of Fresh Readings | 96 |
| 3.4.4 Secrecy of the Template Sketch | 97 |
| 3.5 Secure Remote Fuzzy Extractors from Remote Secure Sketches | 98 |
| 3.6 Conclusion and Future Work | 99 |
| Conclusion and Future Work | 101 |
| Conclusion | 101 |
| Future Work | 102 |
| Résumé par chapitre | 105 |
| Chapitre 2 | 105 |
| Chapitre 3 | 107 |
| References | 111 |
| List of Figures | 119 |
| List of Tables | 120 |

Introduction

Outline of the current chapter

| | |
|---|-----------|
| 1.1 About Biometrics | 2 |
| 1.1.1 History of Biometric | 2 |
| 1.1.2 Biometric Definition | 4 |
| 1.2 About the Data | 4 |
| 1.2.1 Personal Data | 4 |
| 1.2.2 Biometric Data | 5 |
| 1.2.3 Comparison Between Personal Data and Biometric Data | 5 |
| 1.3 Biometric Authentication Systems | 5 |
| 1.3.1 Specification | 6 |
| 1.3.2 Accuracy Metrics | 7 |
| 1.3.3 The Threat Points of a Biometric System | 9 |
| 1.3.4 Typology of Biometric Attacks | 10 |
| 1.4 Cryptographic Methods for Fuzzy Data | 11 |
| 1.5 Biometric System Problems Addressed | 13 |
| 1.5.1 Exhaustive Search Attacks Complexity | 14 |
| 1.5.2 Design of a Remote Secure Sketch Resistant to Offline Attacks | 14 |
| 1.6 Publications, and Submitted Works | 14 |
| 1.6.1 Attacks on Biometric Systems | 14 |
| 1.6.2 Construction of Cryptographic Primitive | 15 |
| 1.6.3 Related Publications | 15 |

"Bond, your fingerprint will unlock the device. Handle it with care." This sentence, pronounced by Q in "Die Another Day" (2002), underscores the integration of biometric technologies as a secure method for accessing sensitive intelligence. This representation mimics their real-world application, where biometrics are relied upon for their perceived reliability and efficacy in safeguarding confidential information and controlling access to critical resources.

Biometric technology is transforming how we secure our digital and physical information by

relying on physical and behavioral characteristics, such as fingerprints, facial features, keystrokes, and voice patterns. In the collective unconscious, biometrics is perceived as a highly sophisticated and extremely secure method of authentication. As highlighted above, this perception is often reinforced by numerous myths propagated by popular culture, including films that romanticize the technique, showcasing how fingerprints, retinal scans, and voice recognition are used to unlock everything from personal devices to secure facilities. These systems become more widespread, they promise to transform industries ranging from finance and healthcare to law enforcement and personal device security. Banks are implementing facial recognition for secure transactions, hospitals are utilizing biometrics to rapidly identify patients and access medical records, and law enforcement agencies are employing advanced biometric databases to enhance public safety. More casually, we can daily observe this technology in action, for example when we unlock our smartphones with a fingerprint or face scan or when we pay for our groceries [5] with our palm vein.

While biometrics systems claim to offer a higher level of security compared to traditional passwords and PINs, is it really the case? To illustrate this point, we may consider the example of Apple's Touch ID [4] to compare the security of a traditional password with the security of biometric technologies. Apple asserts that the False Match Rate (FMR) of Touch ID is 1 in 50,000. In other words, the probability that another individual could unlock your device with their fingerprint is 0.002%. Upon initial examination, this appears to be a highly secure system, but it is less secure than a five-digit password, which has a 1 in 100,000 probability of being guessed *i.e.*, the probability that another individual could unlock your device is 0.001% or approximately 16 bits of security. The problem is that a 5 digit password is not secure. Even for a human who takes 3 second to test a password, it would take 3 days and 12 hours. For a mainstream computer, it would only take 1 to 2 seconds.

Furthermore, biometric systems present a distinct set of challenges and vulnerabilities. Unlike passwords, biometric traits cannot be changed if they are compromised. If a password is leaked, it can be reset; however, if a fingerprint is cloned or stolen, it cannot be replaced. The unchanging nature of biometric data raises concerns about privacy and security. Additionally, biometric systems are susceptible to spoofing attacks, where artificial representations of biometric traits, such as fake fingerprints or masks, are used to deceive the system. The advent of sophisticated technologies, like DeepFake, has exacerbated these risks, enabling more advanced and convincing spoofing methods.

Although this thesis does not address the problem of advanced spoofing, its objective is to identify the inherent vulnerabilities in biometric systems and propose solutions to enhance their security. The primary goal is to provide insights and recommendations that help mitigate the risks associated with biometric authentication in existing systems. Furthermore, this work aims to propose novel cryptographic primitives designed to securely handle biometric authentication, thereby supplementing current methods and paving the way for a more secure future.

1.1 About Biometrics

In this section, we introduce the concept of biometrics, its definition and its history since –5000 BCE.

1.1.1 History of Biometric

The origins of biometrics can be traced back to ancient times when fingerprints were utilized for identification purposes in ancient societies such as China (7th century BCE) and Babylon

(-5000 BCE). From these early beginnings, the development of modern biometric techniques can be summarized in five significant periods:

Origins of Research (1820-1870): During this period, the foundational research into biometrics began to take shape:

- **1820:** Marcello Malpighi's groundbreaking thesis on fingerprints is published, laying the groundwork for future studies in biometric identification.
- **1823:** Johan Evangelista Purkinje demonstrates the unique nature of fingerprints, recognizing their potential for individual identification.

Formalization of Biometrics (1870-1936): The formalization of biometric principles and methods emerged during this period:

- **1858:** Thumbprints are utilized as proof of identity in India by Sir Francis Galton, marking an early application of biometric principles in real-world contexts.
- **1870:** Henry Faulds proposes the use of fingerprints for identification via a publication in the journal 'Nature'.
- **1880:** The Bertillonage method is invented by Alphonse Bertillon, providing a systematic approach to anthropometric identification.

Emergence of Tools and Techniques (1936-1996): In the mid to late 20th century, advancements in biometric technologies began to accelerate:

- **1936:** Frank Burch proposes the use of iris patterns for biometric recognition, introducing a novel approach to identification.
- **1960:** Professor Gunnar explores voice recognition as a potential method for identification, laying the groundwork for future developments in voice biometrics.
- **1970:** The development of automatic fingerprint processing terminals revolutionizes biometric identification, enabling rapid and accurate fingerprint analysis.
- **1980:** Commercial introduction of automatic fingerprint processing terminals further expands the use of biometric identification systems in various industries.
- **1992:** In London, Edward Henry develops the fingerprint identification system, known as the Henry system.

Integration into Mainstream Technologies (1996-2010): Biometric technologies begin to integrate into mainstream technologies:

- **1990:** The release of Dragon Dictate, the first public voice recognition tool, marks a significant milestone in the adoption of voice biometrics.
- **1996:** Hand geometry systems are used for access control at the Atlanta Olympic Games, showcasing the reliability of biometric systems in high-security environments.
- **2010:** Widespread integration of biometric authentication into smartphones enhances user security and convenience.

Recent Advances (2010-2024): In the past decade, biometric technologies have seen widespread adoption and continued innovation:

- **2010:** The resurgence of voice recognition, coupled with the emergence of vocal assistants.

- **2014:** Facebook introduces DeepFace, an algorithm capable of identifying individuals in photos with 97% accuracy, showcasing advancements in facial recognition technology.
- **2020:** Introduction of biometric payment cards replaces traditional PIN codes with fingerprint authentication.
- **2024:** The introduction of biometric payment terminals replaces all payment interactions with a scan of the palm vein.

Over time, significant technological advancements, notably in the fields of computing and imaging, facilitated the development of new techniques and improved accuracy. For further details, refer to these two blog articles [8, 24].

1.1.2 Biometric Definition

Biometrics encompasses a range of computer techniques designed to automatically recognize individuals based on their physical, biological, or even behavioral characteristics [9]. Biometric data enables the identification of an individual, with many of these characteristics being both unique and permanent. The purpose of biometrics is to uniquely identify individuals, thereby serving to protect sensitive data or sites, verify ownership of specific accounts or objects, and mitigate identity theft due to the uniqueness and difficulty of replicating these data. In criminology, biometrics is utilized to differentiate suspects in a case or identify a victim. The scope of applications for biometrics is extensive and expected to continue expanding. Nowadays, biometrics recognition can use various modalities, such as:

- | | | |
|-------------------|-------------|----------|
| • Fingerprint | • Iris | • Retina |
| • Finger geometry | • Face | • Voice |
| • Palm Vein | • DNA | • Gait |
| • Palm geometry | • Keystroke | • Dental |

1.2 About the Data

In this section, we recall what is personal data. We explain how biometric data qualifies as personal data, highlighting its specific characteristics compared to other types of personal data.

1.2.1 Personal Data

Personal data are a range of data defining a person. These data can be used to identify a physical person or to distinguish an individual within a group. This data alone may allow the direct identification of an individual. This is the case, for example, with the first and last names. In some cases, identification can be done indirectly, for example, with the phone number, license plate, or social security number. More generally, any information concerning a physical person is considered personal data. Even the most generic information such as date of birth, city of birth, or gender can, through cross-referencing, accurately identify a physical person. The use of such strategies allows for the potential de-anonymization of databases, using publicly available information or auxiliary data to find the individuals associated with the data as shown in [87, 115, 88]. It is important to specify that personal data only concerns physical entities such as human beings and not moral entities such as companies or associations. In France, the processing of this data is supervised by

the Commission Nationale de l'Informatique et des Libertés also known as French Data Protection Authority, or CNIL.

1.2.2 Biometric Data

Biometric data can be used to (uniquely) identify a person. As a result, these data are personal data that possess all of the aforementioned characteristics and additional nuances. For instance, biometric data are:

- *Sensitive*: Biometrics are generally stable and do not change significantly over time, providing a reliable means of identification over time. Biometrics can reveal more personal information than expected. For example, some medical conditions can be identified using biometric information, as shown by Ross *et. al.* [32]. Other less significant information, known as soft biometrics, can be disclosed, such as gender, hair color, iris color, height, or even handedness [47], depending on the used biometric system.
- *Used for security purpose*: Because of their stability, these data are used to authenticate individuals to various systems, such as banks. In most cases, biometric authentication is not sufficient to provide an adequate level of security, which requires the use of a second factor. One of the problems with biometrics is that with the popularity of social media, some of our modalities are public, such as face or voice. Database leaks can also expose non-public biometric modalities such as fingerprints to a government database. With the development of new techniques such as DeepFake [52, 26], the risk of our public biometrics being used to impersonate us is greater than ever.
- *Variable*: Although biometric data demonstrate stability, they differ from passwords in their inherent variability. This characteristic renders conventional cryptographic techniques, such as hash functions, inadequate. Consequently, novel methodologies tailored to the unique nature of biometric data must be developed.

1.2.3 Comparison Between Personal Data and Biometric Data

To emphasize the position of biometric data within the broader category of personal data, we present a comparative analysis in Table 1.1. It provides a comprehensive overview of the characteristics that distinguish biometric data from other types of personal data. While both types of data can be used to identify individuals, biometric data offers a more direct and unique means of identification due to their stability and specificity. However, this also renders biometric data more sensitive and, as a subset of personal data, they require specific security and privacy regulations.

1.3 Biometric Authentication Systems

In this section, we recall what constitutes a biometric authentication or identification system. We outline its specifications, including the properties it must comply with and how its reliability is assessed by detailing accuracy metrics such as the False Match Rate (FMR). Finally, we explore the various methods for attacking a biometric system and provide a categorization of the different types of attacks.

| Characteristic | Personal Data | Biometric Data |
|---------------------------------|--|---|
| Definition | Information relating to an identified or identifiable person | Unique physiological or behavioral characteristics that can identify an individual |
| Examples | Name, phone number, social security number, date of birth | Fingerprints, face, iris, voice patterns |
| Identifiability | Can directly or indirectly identify a person | Can directly identify a person with high accuracy |
| Sensitivity | Varies; some data are more sensitive (<i>e.g.</i> , social security number) | Highly sensitive; can reveal medical information |
| Stability | Varies; some data (<i>e.g.</i> , phone number) may change | Generally stable over time (<i>e.g.</i> , fingerprints) |
| Security Usage | Used for identification and verification | Used for authentication, identification and access control |
| Privacy Risk | Can be anonymized or pseudonymized | Can be anonymized or pseudonymized |
| Public Exposure | Generally private; can be protected with encryption and access controls | Some modalities (<i>e.g.</i> , facial features) can be publicly available; higher risk of exposure |
| Regulation | Subject to data protection regulations (<i>e.g.</i> , GDPR, CNIL) | Subject to data protection regulations (<i>e.g.</i> , GDPR, CNIL) |
| Cryptographic Techniques | Can use traditional cryptographic methods (<i>e.g.</i> , hashing, encryption) | Requires specialized techniques due to inherent variability |

Table 1.1 – Comparison Between Personal Data and Biometric Data

1.3.1 Specification

In a biometric recognition system, biometric templates (*i.e.*, a transformed version of their biometric information) of users are stored in a database. The first operational mode (identification) involves determining the identity of an individual by comparing their freshly provided template with all the templates stored in the database. The second mode, authentication (verification), corresponds to verifying the claimed identity by comparing the corresponding enrolled template with the fresh template provided by the user. In both cases, the enrollment process remains the same. The sensor captures the image of the biometric modality. This measurement then undergoes a sequence of transformations, including feature extraction (*e.g.*, using Gabor filtering [140, 134]) followed by a Scale-then-Round process [38] to convert the data into a format better suited for cryptographic schemes, such as binary or integer vectors. These templates are then protected either through encryption alone or by utilizing Biometric Template Protection (BTP). Figure 1.1 illustrates the entire process.

Properties for a Biometric Scheme

The purpose of a biometric scheme is to generate a template. Similar to a hash function, the transformation leading to this template must ensure several properties that preserve the privacy of the biometric data. To mitigate the problems caused by a database leak, the essential security and

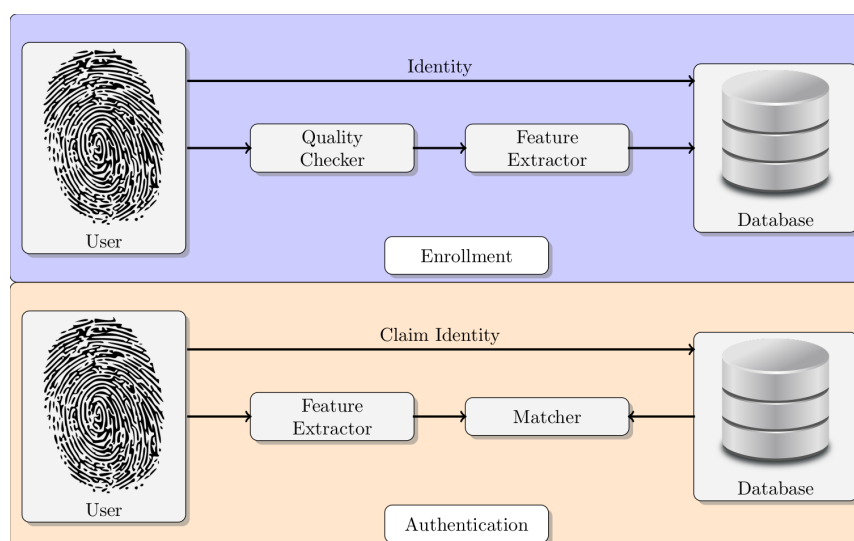


Figure 1.1 – Schematic view of a biometric system for enrollment and authentication.

performance criteria that must be met by biometric recognition systems are identified in ISO/IEC 2474 [1] and ISO/IEC 30136 [2]:

- *Irreversibility*: The transformation process should be irreversible, meaning that it should be impossible to derive the original biometric data from the template. Using the analogy to hash functions, this definition can be extended to pre-images, signifying that for an attacker, it should be difficult to find data such that its transformation yields the targeted template.
- *Unlinkability*: Biometric templates should not reveal any information that can link them to a specific individual, thereby ensuring privacy and anonymity. In a broader sense, when comparing templates from two different services, we should not be able to determine if the templates are derived from the same biometric data. Analogous to hash functions, the templates must be indistinguishable from random strings.
- *Revocability*: It should be possible to remove a stolen template and substitute it with a new one. This measure serves to prevent the illicit use of the stolen template for unauthorized access to protected resources while enabling the legitimate user to maintain uninterrupted access. In practice, without the usage of a revocable second factor (*e.g.*, password), this property is difficult to fulfill.
- *Performance preservation*: The transformation process should maintain the performance of the biometric system, ensuring accurate recognition and verification.

1.3.2 Accuracy Metrics

To assess the security of a biometric system, different metrics are used based on the operation mode (recognition or identification). In this section, we recall the most commonly used metrics to assess the accuracy¹ of biometric systems in a non-exhaustive manner.

In the context of authentication systems (specifically, in a 1:1 system), the predominant metric utilized is the False Match Rate (FMR). This rate serves as an empirical estimation denoting the likelihood of a biometric sample being incorrectly recognized by the matcher. In other words,

1. In the field of biometrics, the term "accuracy" is employed to describe the probability that a recognition algorithm may incorrectly identify a person when attempting to identify or authenticate them. The term should therefore be distinguished from its meaning in machine learning.

this is an estimation of the probability that the matcher incorrectly decides that a newly collected template matches the stored reference. According to the Face Recognition Technology Evaluation (FRTE) 1:1 Verification [14], given a vector of N imposter scores v and T a threshold, the FMR is

$$\text{FMR} = \frac{1}{N} \sum_{i=1}^N S(T - v_i)$$

with $S(\cdot)$ the unit step function, and $S(0)$ taken to be 1. Similarly, given a vector of N genuine scores u , the False Non Match Rate (FNMR) is computed as,

$$\text{FNMR} = 1 - \frac{1}{N} \sum_{i=1}^N S(u_i - T).$$

In other words, the FNMR gives an estimation of the probability that a genuine sample is incorrectly rejected by the matcher.

The confusion between the False Match Rate (FMR) and the False Acceptance Rate (FAR) often arises due to their subtle distinction. The FAR operates at the system-wide scale, encompassing more than the matching component. Specifically, it represents the probability of a biometric sample being falsely recognized by the entire system. The FAR considers the collective performance of all security layers within a biometric system, such as liveness detection, thereby providing a comprehensive assessment of system integrity. Let FTA denote the Failure To Acquire rate, *i.e.*, the probability that the system fails to produce a sample of sufficient quality. We have the following equality

$$\text{FAR} = \text{FMR} \times (1 - \text{FTA}).$$

The FTA is composed of different error tests *e.g.*, the Failure to Capture (FTC) *i.e.*, the sensor cannot successfully detect a sample and the Failure to Extract (FTX) *i.e.*, the sample's quality is not good enough to generate a valid template.

The misconception between the False Non-Match Rate (FNMR) and the False Reject Rate (FRR) persists for analogous reasons as discussed above. Those two notions can be linked by

$$\text{FRR} = \text{FTA} + \text{FNMR} \times (1 - \text{FTA}).$$

The metrics mentioned earlier rely on a threshold selected to minimize either the False Non-Match Rate (FNMR) or the False Match Rate (FMR). Typically, this threshold is determined at the Equal Error Rate (EER) where the FNMR equals the FMR.

In the context of identification (specifically, in a 1: N system), a widely used metric is the False Positive Identification Rate (FPIR). This metric quantifies the error rate when the system misidentifies an impostor as a user. The False Non-Identification Rate (FNIR) assesses the likelihood of genuine users being incorrectly rejected or failing to be identified by the system in the identification mode (1: N). It is important to note that the two definitions above (FPIR and FNIR) hold for a threshold identification system, *i.e.*, a system that tries to recover the identity of the client in the database or rejects it if no one known is close enough. These definitions can be relaxed for a system that does not reject any user, but solely returns the closest one in the database.

More errors and metrics exist, as shown in the NIST report [34] and this blog article [25]. However, as these measures are not relevant to this thesis, we do not intend to develop them

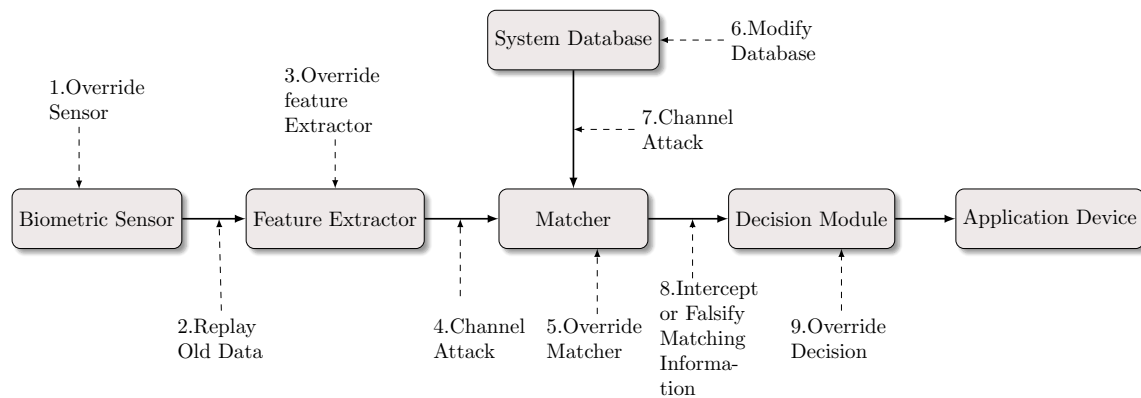


Figure 1.2 – Attack points in a generic biometric recognition system.

further.

1.3.3 The Threat Points of a Biometric System

In the context of a biometric system, a critical vulnerability arises when information is intercepted between the matcher and the decision module, as illustrated in Figure 1.2 (Point 8). This figure, based on Ratha *et al.*'s [133] research, introduces both the decision module and two additional threat points. Note that for enhanced readability, the Pre-Processing, Feature Extraction, and Template Generation have been grouped into the Feature Extractor. The control of Point 4 allows the submission of a chosen template, while Point 8 grants access to additional information beyond the binary output yes/no. Each point can be used for different purposes, for instance:

- Point 1 is the most used because it requires the least knowledge about the attacked system. It allows the attacker to perform presentation attacks [50] and to present a potential replica of the biometric data to the system. such as a fake finger [131, 130], fake palm vein [78] a forged signature [56], a fake iris [114, 112, 76], a facial mask or a morphed image [69, 82, 48]. Moreover, this point allows the attacker to bypass all the sensor security systems such as liveness detection [73, 118] or sensing finger conductivity or pulse [27].
- From Point 2, an attacker can either read the biometric data if the system is in use or replay an old biometric sample. The sensor could detect an old biometric sample, but at this point, the sensor is bypassed.
- With the control of Point 3, an attacker can override the extractor and generate the feature sets or templates of its choice.
- As specified earlier, Point 4 allows the capability to either submit any feature vector or template, or to read the template or feature vector of the current user. The two steps of feature extraction (with or without transformation) and matching are often inseparable, making this mode of attack highly challenging. However, when minutiae are transmitted to a remote matcher, the threat becomes significant. In such cases, an attacker could potentially intercept the communication and modify certain packets.
- For Point 5, the matcher is manipulated to consistently produce the desired outcome for the attacker, whether it be high or low match scores.
- With control of Point 6, the attacker can manipulate the database. This includes adding, removing, or tampering with a template of their choice, potentially leading to unauthorized access for an impostor, or at the very least, denial of service for a legitimate user.

- When the stored template or feature vector is sent to the matcher, the control of Point 7 allows the attacker to either read it or modify it. The issues of such control are the same as Point 4 described above.
- By exploiting Point 8, the attacker can intercept matching information, such as the score, or modify it to their advantage.
- Point 9 is one of the most critical threat points. If the outcome can be manipulated by the attacker's chosen result, it could lead to critical consequences. Even if the system itself shows excellent performance, its effectiveness is undermined by the ability to override the outcome.
- An additional threat point could be added for the communication between the decision module and the application. However, the impact of such control is the same as Point 9, and thus it is not further discussed.

It is natural to assume that the attacker controls only a single point, but nothing is preventing it from controlling several of them. Some points are more sensitive, while others are more reachable by the attacker. For more detailed insights into the threat points, readers are referred to Ratha *et al.*'s work [133]. Note that in Chapter 2, we suppose that the attacker is in control of Point 3 or 4 in the sense that it may submit the template of its choice. The fact that the attacker can submit the template of its choice allows us to directly evaluate the FMR of the system rather than the FAR, making these two points a very interesting case study as the FMR assesses directly the accuracy of the matcher without any additional defense features (see Section 1.3.2).

1.3.4 Typology of Biometric Attacks

To compile the attacks on biometric systems, several state-of-the-art exist. For example, "*Attack and Presentation Attack Detection*" [50] introduces attacks and their detection with more than a hundred sources. Furthermore, for cancelable biometrics (a widely used biometric), "*Cancelable Biometrics*" [77] presents the different existing methods as well as a small panel of attacks. For more specific modalities, "*Face Recognition Systems Under Morphing Attacks: A Survey*" [51] gives the state of the art of spoofing in face recognition attacks and some countermeasures. It exists also the same work for fingerprint in "*Security and Accuracy of Fingerprint-Based Biometrics: A Review*" [53] where the authors present two attacks to biometric systems and countermeasures. For cryptography and watermarking, "*On the Vulnerability of Biometric Security Systems*" [124] presents an overview of the weakness of biometric security systems and possible solutions to improve it. All the attacks presented in the state of the art above or any future attacks can be categorized according to the threat levels described by Simoens *et al.* [91]. The threat levels are categorized as follows:

1. *Biometric reference recovery*: The adversary aims to retrieve the stored biometric template. In this category, among others, we find the center attack [83] and our work on the leakage exploit [19] (see Chapter 2).
2. *Biometric sample recovery*: The adversary aims to generate a new biometric template that successfully passes authentication within the biometric system. In this category, among others, we find morphing attacks [51], spoofing attacks, presentation attacks [50], and exhaustive searches.
3. *Tracking users with different identities*: This type of attack can occur when various references from the same user, potentially originating from different applications, can be correlated. A system that can withstand such attacks is considered to offer identity privacy [113].

4. *Tracking users over different queries*: The adversary seeks to link queries. The characteristic of a system that thwarts such an attack is referred to as transaction anonymity [113].

Note that each of those threat levels echoes a violation of an ISO property described above (see Section 1.3.1). The methodologies employed in the execution of these attacks may be classified into five distinct categories:

- Hardware Based
- Data Based
- Template and Matcher Based
- Side Channel
- Server level

Figure 1.3 from Sharma *et al.* [23] gives an insight on those categories.

1.4 Cryptographic Methods for Fuzzy Data

Biometric Template Protection is a scheme aiming to protect biometric features by transforming them into protected templates. These functions behave similarly to hash functions even if close verification must be achievable and respect the properties of ISO (see Section 1.3.1) The main schemes of BTP can be categorized into three distinct approaches:

- *Biometric Cryptosystems* (BC).
- *Cancelable Biometrics* (CB).
- *Keyed Biometrics* (KB).

Cancelable Biometric: This approach modifies the original biometric representation in an irreversible and controlled manner, preferably using irreversible transformations that can be parametrized with public salts or user-specified secrets. Even if biometric data are compromised, they cannot be used for direct authentication. Instead, a new biometric representation is generated from the original data each time authentication is performed. The output of this transformation is stored on the server, allowing verification of the user by comparing the transformed fresh biometric data with the transformed enrolled biometric data. CB employed several methodologies to perform the transformation, thereby enabling Manisha and Nitin [42] to create a recent taxonomy of CB techniques. In this taxonomy, there are six categories, such as:

- **Cryptography based** including among others BioHashing, Steganography, Index of Max Hashing, and Fuzzy Commitment.
- **Transformation based** including among others Random Projection, Wavelet Transformation, and Rotation.
- **Filter based** including among others Bloom Filter and Gabor Filter.
- **Hybrid Based** including the Template Transformation followed by a Biohashing.
- **Multimodal based** including among others Random Projection followed by a transformation and Bloom Filter followed by a feature level fusion.
- **Other methods** where there are among others Physically Unclonable Functions, Binary Gaussian Mixture, Huffman Encoding, and Deep Learning.

The taxonomy in question suffers from some flaws, but to cite one, steganography is not a subset of cryptography. Rather, it is a distinct field that employs completely different methodologies and has completely different applications (see Table 1.2). Because of these problems, we prefer the taxonomy proposed by Patel *et al.* [77] which has the following categories:

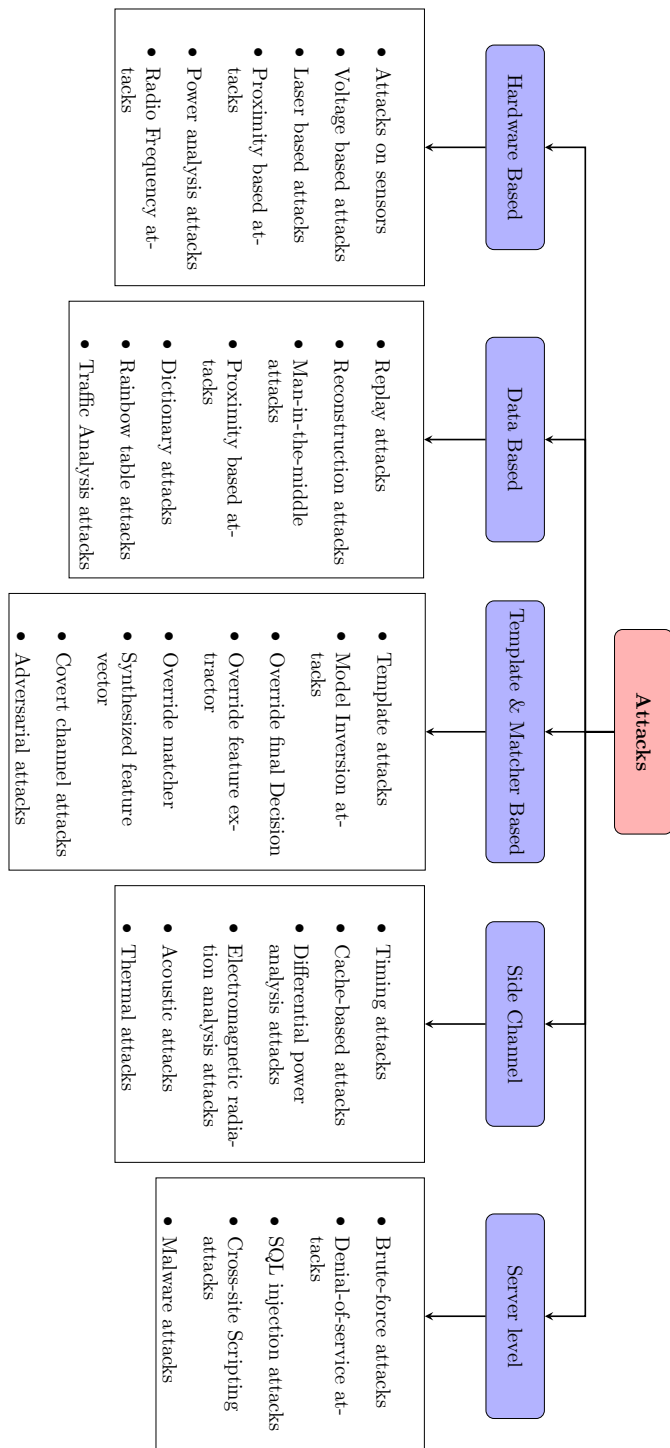


Figure 1.3 – Classification of attacks on biometric cryptosystems (Sharma *et al.* [23]).

| | Steganography | Cryptography |
|---------------------------|---|---|
| Definition | Technique to hide information within another object | Technique to convert plain text into ciphertext |
| Objective | Hide the message existence | Protect the message meaning |
| Visibility of the message | Never | Always |
| Failure | Message may be discovered | Message gets decrypted |

Table 1.2 – Short comparison between Steganography and Cryptography

- **Cancelable Biometric Filters**
- **Knowledge Signatures**
- **Random Permutations**
- **BioHashing Methods**
- **Random Projections**
- **Bioconvolving**
- **Bloom Filters**
- **Geometric Transforms**
- **Salting Methods**
- **Hybrid Methods**

Concerning the performances and the accuracy, despite the variety of the existing techniques as shown above, they are unsatisfactory. Indeed, according to Patel *et al.* [77], the best performances of CB give a EER close to 2% and FAR of 10^{-4} while other methods yields way better accuracy (e.g., FMR of 10^{-6} for [3]). For more details, the reader is referred to the surveys [77, 42, 23].

Keyed Biometrics: In KB, secure computation techniques are used to enable the verification in the encrypted domain. The incorporation of a cryptographic key into the biometric data serves to enhance the security of the system. To authenticate, both the correct biometric feature and the corresponding key must be present. This introduces an additional layer of security by incorporating a secret element into the authentication process. This category includes, for example, Jarrous and Pinkas [105] and Bringer *et al.* [80]. To the best of our knowledge, there is no state-of-the-art, survey or SoK on Keyed Biometrics which may be the subject of future research.

Biometric Cryptosystems: BC represents a biometric security approach that integrates cryptographic techniques to guarantee the security and confidentiality of biometric data. In BC, a cryptographic key K is bound to a biometric input x through a probabilistic algorithm that takes a randomizer r as additional input. A *helper data* HD is derived from x and is stored either in a remote database or on an end-user device. This helper data serves as a protected template and should not give any information on x . BC schemes include fuzzy commitments [135], fuzzy vaults [116], secure sketches [122] and fuzzy extractors [123, 17]. Those primitives serve two different purposes. Fuzzy commitments, fuzzy vaults, and fuzzy extractors aim at reproducing the key K using another biometric input $x' \approx x$ and r , HD. The secure sketch is used to recover x using another biometric input $x' \approx x$ and HD. Secure sketches and fuzzy extractors are further investigated in Chapter 3. For more details, the reader is referred to the survey [23].

1.5 Biometric System Problems Addressed

This thesis examines exhaustive search attacks in depth. Given that these attacks are predominantly offline rather than online due to online countermeasures (e.g., limiting the number of

attempts that a password can be tried or introducing time delays between successive attempts), to mitigate those attacks, we present a protocol resistant to offline attacks. This section provides a brief introduction to the work presented in the subsequent chapters.

1.5.1 Exhaustive Search Attacks Complexity

In Chapter 2, we analyze the security and efficiency of biometric systems through several axes. We consider two setups. In the first, the attacker has access to the database in clear text. In this case, we show how he can make a partition of the database to reduce the number of templates he has to reverse to impersonate all the users. This configuration allows us to investigate the probability of near-collisions. In the second, the database is either on a server or in a trusted environment on a device. In this case, we evaluate the number of tests an attacker must perform to impersonate a user. We then analyze the security of a template database using the same methodology as we would for a password database. We provide an analysis of untargeted attacks and collisions.

Remark 1.5.1. *A major part of those results are from Durbet et al. [30, 18, 19]. The remainder is on going research in collaboration with Paul-Marie Grollemund, Pascal Lafourcade and Kevin Thiry-Atighehchi.*

1.5.2 Design of a Remote Secure Sketch Resistant to Offline Attacks

In Chapter 3, we propose several protocols along with a computational secure sketch. Firstly, a novel secure sketch is created with the objective of facilitating the reconstruction of stored data. The security of the stored data is guaranteed by a computational assumption. Subsequently, this construction is employed to generate a zero-knowledge protocol. We propose an authentication protocol that allows for the authentication of a service without disclosing any information other than the client's knowledge of the secret. We then propose a protocol that enables the secure correction of fresh data in collaboration with a remote server. Under computational hypothesis, we guarantee the confidentiality of both the fresh data and the enrolled data. Finally, we demonstrate the utilization of our constructs to develop generic fuzzy extractors, thereby enabling the secure reproduction of a secret with a remote server under the same hypotheses.

Remark 1.5.2. *A major part of those results are from Durbet et al. [12]. The remainder is on going research in collaboration with Koray Karabina and Kevin Thiry-Atighehchi.*

1.6 Publications, and Submitted Works

In this section, we present our contributions through published papers, submitted papers, and work in progress. For each work, we provide the abstract and detailed contributions, organized thematically for clarity.

1.6.1 Attacks on Biometric Systems

In this section, we present our contribution and works revolving around generic attacks in biometrics.

Near-collisions and Their Impact on Biometric Security (Published): This paper [30] was published in the Proceedings of the 19th International Conference on Security and Cryptography (SECRYPT) and presented in Lisbon, Portugal. This work has been done in collaboration with Paul-Marie Grollemund, Pascal Lafourcade and Kevin Thiry-Atighehchi. It extends our previous work by examining the mechanisms enabling the impersonation of multiple individuals and the impact of database size on near-collisions. We developed an efficient algorithm for partitioning a leaked template database, enhancing the generation of a master-template set. The practical implications of our algorithms are demonstrated through experimental studies. This work is detailed further in Section 2.2 of this manuscript.

Biometric Untargeted Attacks: A Case Study on Near-collisions (Submitted): This work [18] is under peer review. This work has been done in collaboration with Paul-Marie Grollemund and Kevin Thiry-Atighehchi. It advances our research on near-collisions and exhaustive search attacks by focusing on untargeted attacks. We use probabilistic modeling to assess the theoretical security limits of biometric systems, considering metric space and system parameters. This work is detailed further in Section 2.5 of this manuscript.

Exploit the Leak: Understanding Risks in Biometric Matchers (Submitted): This work [19] is under peer review. This work has been done in collaboration with Paul-Marie Grollemund, Dorine Chagnon and Kevin Thiry-Atighehchi. It focuses on exhaustive search attacks on biometric matchers, exploiting various types of information leakage, such as distance values. We present a catalog of information leakage scenarios and their impacts on data privacy, leading to unique attacks with quantified computational costs. This work is detailed further in this manuscript in Section 2.4.3.

1.6.2 Construction of Cryptographic Primitive

To address the security issues identified in the aforementioned section, we present the development of novel cryptographic primitive.

Generic Construction of Secure Sketches from Groups (Submitted): This work [12] is currently under peer review. This work has been done in collaboration with Koray Karabina and Kevin Thiry-Atighehchi. We developed a secure method for using biometric data by introducing a family of secure sketches constructed generically from groups with unique factorization properties. We provide a detailed analysis of the mathematical structures and establish several computational and decisional hardness assumptions. Our secure sketches are efficient, handle a linear fraction of errors with respect to the L_1 distance, and are both reusable and irreversible. This work is detailed further in this manuscript in Chapter 3.

1.6.3 Related Publications

In this section, we present a publication that is not included in this thesis, but which we believe to be of interest.

Authentication Attacks on Projection-based Cancelable Biometric Schemes (Published): This paper [29], published in the Proceedings of the 19th International Conference on Security and

Cryptography (SECRYPT), was presented in Lisbon, Portugal. This work has been done in collaboration with Paul-Marie Grollemund, Pascal Lafourcade, Denis Migdal and Kevin Thiry-Atighehchi. We formalize an attack that reverses a biohash template back to the original biometric image using integer linear programming (ILP) and quadratically constrained quadratic programming (QCQP). This allows an adversary to alter their fingerprint image to impersonate any individual and, in severe cases, simultaneously impersonate multiple individuals. This concept of simultaneous impersonation is fundamental to many of our other works.

The Security of Biometrics Systems Against Targeted Attacks, Untargeted Attacks and Near-Collisions

Outline of the current chapter

| | |
|--|-----------|
| 2.1 Introduction | 18 |
| 2.2 Master Templates and Near-Collisions | 21 |
| 2.2.1 Definitions | 21 |
| 2.2.2 Database Partitioning | 22 |
| 2.2.3 Multi-Near-Collisions and Master-templates | 31 |
| 2.2.4 Risk related to Master-templates | 36 |
| 2.2.5 Numerical Evaluations: Databases Security w.r.t Near-Collisions . . . | 38 |
| 2.2.6 Metric Based Analysis | 38 |
| 2.3 Targeted and Untargeted Attack Models | 42 |
| 2.4 Targeted Attack | 43 |
| 2.4.1 Exhaustive Search | 43 |
| 2.4.2 The Center Search Attack | 45 |
| 2.4.3 The Impact of Information Leakage on Targeted Attack | 46 |
| 2.4.4 Exploiting the Leakage | 48 |
| 2.4.5 Accumulation Attack: A Passive Attack | 50 |
| 2.4.6 Weaker Attack Model: Compromise a point in the ball | 52 |
| 2.5 Untargeted Attack | 52 |
| 2.5.1 Metric Space Based Bounds | 52 |
| 2.5.2 Biometric security: Assess the Security Score Against the Outsider . . | 59 |
| 2.6 Conclusion and Future Work | 60 |

Abstract of the current chapter:

In this chapter, we examine the security of biometric data against various attacks. Biometric data, after transformation, behave similarly to passwords. Therefore, we treat them with the same approach as passwords. Password analysis revolves around three main axes, which we apply here.

The first aspect is collision search *i.e.*, finding two distinct biometric data from two distinct users leading to the same template. A major difference between biometric data and passwords, is the intrinsic unpredictable variability of biometric data, even for the same individual. As a result, in the case of biometric data, the system must tolerate variations, and exact matches are not relevant, which is standard for passwords. Then, instead of exact collisions, we consider near-collisions. We study and characterize near-collisions, which are problematic for both the performance and security of biometric systems. A high number of near-collisions significantly degrades system accuracy. In addition, we show how an attacker can exploit near-collisions to create templates capable of impersonating multiple users. To mitigate this issue, we determine their occurrence probabilities and derive a scoring method to aid in biometric system parameterization.

The second aspect is targeted exhaustive search attacks, which focus on a specific user. Although this work has been initiated in the literature, we extend it to account for different attack scenarios. We characterize various data leakage scenarios, such as distance leakage or error position leakage. Each scenario leads to generic attacks for which we provide complexity analyses. This enables the design of biometric systems with a clear understanding of the actual complexity of exhaustive search attacks based on the matcher's information leakage.

The final aspect is untargeted exhaustive search attacks, which do not aim at any specific user. This is crucial because in a highly populated database, even if the probability of impersonating a specific individual remains low, the probability of impersonating someone increases significantly. We characterize the probability of an individual in the database having their identity impersonated. This characterization allows us to derive a security score against such attacks and provide recommendations to limit their impact.

2.1 Introduction

Context: Biometric systems may be divided into two main categories. The first one utilizes plain-text data to facilitate a decision. These types of systems are not designed to preserve user privacy, but they are (highly) accurate. In this chapter, we focus on the second type of system, where the biometric data is not in plain text. Although these systems are less precise than their counterparts, they nevertheless ensure user privacy. In order to achieve this, these algorithms employ a transformation function that takes the biometric data as input, with or without any additional information and outputs a protected biometric data called a template. These transformations must adhere to the ISO [1, 2] specifications. The specifications include four essential properties. *Irreversibility:* Given a template, it must be difficult to retrieve the biometric data. This property can be extended to mean that, given a template, it must be difficult to find a biometric data item which, through the same transformation, gives the same or a close template with respect to the system threshold. *Unlinkability:* Given two templates, it must be difficult to tell whether they come from the same user. *Revocability:* It refers to the ability to revoke or deactivate a compromised template and

generate a new one. The final property is *performance preservation*, which stipulates that a transformation should not result in a significant loss of accuracy. The objective of this transformation is to utilize biometric data as passwords while maintaining their privacy (See Section 1.3.1 for more details). The protection of passwords is ensured by the use of hash functions that respect two major properties. The *irreversibility* and *collision resistance* properties. These two properties are identical to those of template *irreversibility*, except that for hashes, the attacker seeks perfect matches. Moreover, hashes can easily be revoked by modifying the salt or the password. Given their analogous properties and applications, it is reasonable to *de facto* consider templates to be akin to hashes. In the literature, there are three principal approaches to the study of hash functions. The first one is to study the occurrence of collisions. The second approach examines the complexity of an attack on a specific user *i.e.*, targeted attacks. The third approach examines the security of the system when all users are simultaneously attacked *i.e.*, untargeted attacks. In this chapter, we apply this methodology to biometric templates to estimate and bound their security concerning near-collisions, targeted, and untargeted attacks.

Detailed Contributions: The contributions presented in the following list provide an overview of the different facets addressed in our work, offering a comprehensive understanding of the results obtained.

- Section 2.2: Our primary contribution is an efficient partitioning algorithm that accelerates attacks aimed at generating a master key or master feature vector. Numerical studies on the implementation of the proposed algorithm demonstrate a reduction in computational time in certain settings. Additionally, we demonstrate a link between the closest string problem with an arbitrary number of words and provide a solution using Simulated Annealing (SANN). Moreover, we determine a bound on the size of a database in function of the template space dimension and the decision threshold, thus preventing near-collisions with a high probability. We introduce the notion of weak near-collisions and strong near-collisions, which enable us to provide a theoretical analysis of the security strength of biometric transformation schemes. The bounds on the probability of a near-collision highlight the theoretical limits on the accuracy of a biometric system. Based on near-collisions analysis, we provide the probability of occurrence of a master template, highlighting the potential for an attacker to impersonate multiple users at once. This analysis has been provided in both setup *i.e.*, targeted and untargeted. We then provide a score that can be used to fine-tuning a biometric system. Finally, we examine works that present near-collisions scenarios using metrics such as FMR, provide the critical population for a given FMR, and provide the minimal FMR for given population sizes.
- Section 2.4: This section presents an analysis of potential information leakage in distance evaluation, with a specific focus on threshold-based obfuscated distance. The contributions include a variety of information leakage scenarios, the corresponding generic attacks, their complexities, and a correction of a result presented in [83]. The aforementioned scenarios give rise to new attack scenarios.
 1. We investigate a novel attack, named accumulation attack, where an *honest-but-curious* server accumulates knowledge during client authentication. This type of attack occurs when there is a minor, yet non-negligible, amount of information leakage.
 2. We introduce new attack strategies by malicious clients that exploit various levels of information leaks from the matcher. Our complexity results, which detail the cost of

these attacks, apply to both *offline exhaustive search attacks* that leverage a leaked (yet obfuscated) database and *online exhaustive search attacks* involving direct interactions with the server.

- Section 2.5: The presented attacks are based on exhaustive search (*i.e.*, brute force attacks) and require only the minimum leakage of information, namely a bit of information about the success of impersonation. Hence, they are possible regardless of the employed BTP scheme, protocol, or biometric modality. We assume a biometric system that makes the best use of the underlying metric space in order to provide theoretical bounds on the complexity of exhaustive search attacks. We use probabilistic modeling to present two matching attack scenarios with the associated security bounds and discuss the security of a template database. The first one, called the “Outsider Scenario”, captures the case where an individual unregistered in a service attempts to impersonate a non-specific user of this service. Specifically, we consider the possibility of an attacker sequentially adapting her strategy. The second scenario, termed the “Multi-Outsider Scenario”, encapsulates cases where several attackers attacks the service in parallel. The bounds on the complexity of the untargeted attacks provide the maximum achievable security. We make recommendations concerning the security parameters during the fine-tuning of a recognition system. Finally, we make recommendations concerning the security parameters during the fine-tuning of a recognition system.

A major part of those results are from Durbet *et al.* [30, 18, 19].

Scope of the results: Our results on the complexity of attacks and the probability of near-collision occurrences apply to many BTP schemes. To the best of our knowledge, many BTP schemes of the three categories (CB, BC and KB) are vulnerable to *offline* attacks regardless of the considered modality. For instance, among the BC schemes, we can identify fuzzy commitments [137], fuzzy vaults [116], and fuzzy extractors [123, 67, 17], to name but a few. Concerning the attacks, they can be performed either *online* or *offline*, and the derived complexity results apply in both cases. Offline attacks are made possible when the protected biometric database is leaked, as the attacker exploits some (even minimal) information that allows her to test a guess. Near-collision yields a theoretical limit on the performance of biometric recognition algorithms.

Intuition of the Mathematical Formulation: After transformation, biometric data are represented as vectors, or templates which are traditionally binary strings. To identify vectors within a maximum distance ε , representing our threshold, we use a ball of radius ε around each template to denote matching points. In other words, if another biometric data point falls within the ball of a given template, the two biometric data match. Thus, a biometric database is a collection of vectors in a metric space, as illustrated in Figure 2.1. Note that this figure is a two-dimensional projection and may not fully capture the true structure. Nonetheless, these representations is used throughout this chapter to help understanding our results and methods. When two templates are close enough that each falls within the other’s ball, they are said to be in near-collision, as depicted in Figure 2.6. Near-collisions are common for the same user as it need to be in near-collision with its template to be authenticated, but if they occur between different users, it may lead to impersonation. If the intersection of their balls is non-empty, a point within this intersection can be found that is close to both templates. This point is referred to as a master template and it can impersonate multiple individuals. Figure 2.2 provides a representation of a master template.

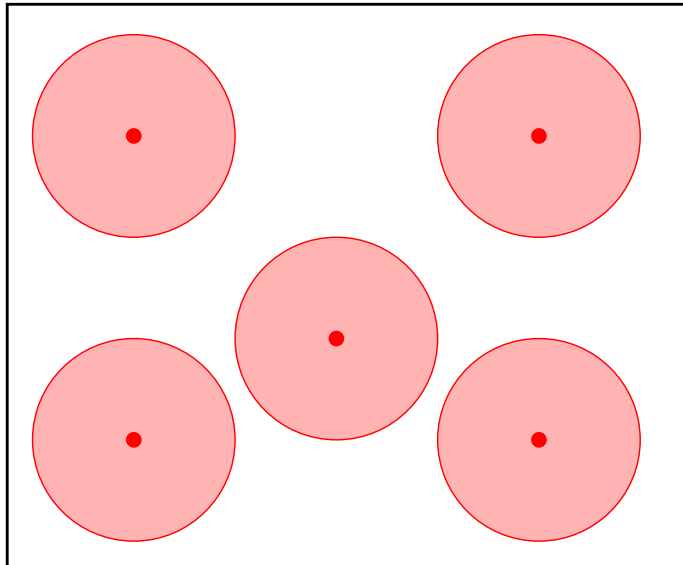


Figure 2.1 – Representation of a database with 5 templates marked by a dot.

2.2 Master Templates and Near-Collisions

In this section, we present a master template for an approach to database partitioning in scenarios where the attacker has access to the database. This may occur when the server is compromised or when a user’s device is stolen. The objective of database partitioning is to reduce the number of templates that an attacker must reverse in order to gain access to the database. In this context, the study of near-collisions emerges as a means of assessing the security of hashed passwords. Near-collisions represent instances where two distinct inputs produce nearly identical templates, which hinders both security and accuracy. We provide bounds on the probability of Near-collisions to allow for the fine-tuning of parameters in a biometric system.

Remark 2.2.1. *This work is analogous to the study by Gernot and Lacharme on the topic of masterkeys [31]. The primary distinction between the two approaches lies in their respective objectives: while their approach seeks feature vectors that generate distinct user templates through bihashing (i.e., construct a biometric feature capable of impersonating multiple users using their password), our approach focuses primarily on biometric templates and does not employ any specific transformation.*

2.2.1 Definitions

As the template space is a metric space, we denote it as (\mathbb{Z}_2^n, d_H) with d_H the Hamming distance.

Definition 2.2.1 (Template database or TDB). *Let (Ω, d) be the template space equipped with the distance d . A subset $\mathcal{B} \subset \Omega$ such that $\mathcal{B} \neq \emptyset$ is a (TDB), or just a database.*

As with hash functions, a preimage to a given template can be found by exhaustive search or other methods depending on the transformation. Then, a natural goal of an attacker is to find a preimage to each template in a leaked *database*. However, as the operation of finding a preimage remains costly, the attacker should focus on ε -master-templates to decrease the complexity of this attack.

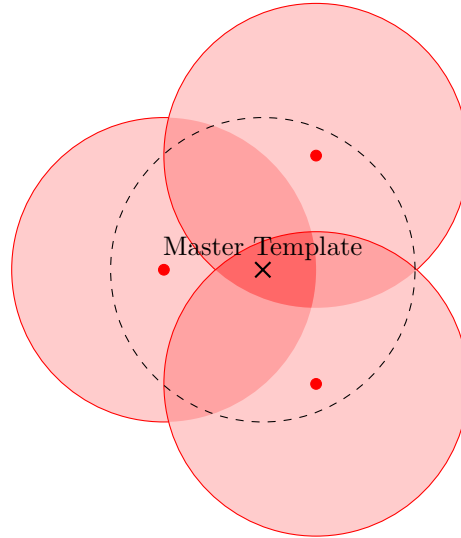


Figure 2.2 – Representation of a ε -master-template marked by a cross for a database with 3 templates marked by a dot.

Definition 2.2.2 (ε -master-template). Let (Ω, d) be the template space and \mathcal{B} be a database. An ε -master-template of \mathcal{B} is $x \in \Omega$ such that $d(x, a) \leq \varepsilon, \forall a \in \mathcal{B}$.

Note that given any TDB, finding an ε -master-template is not always possible. As a counterexample, consider a *database* with two templates such that the distance between them is greater than 2ε . In this case, it is not possible to find a point that is distant by ε from both templates. Figure 2.2 depicts an ε -master-template for a database.

2.2.2 Database Partitioning

This part aims to determine the smallest ε -master-template set for a given *database* \mathcal{B} . To do so, we propose an algorithm based on a clustering algorithm *i.e.*, the Hierarchical Agglomerative Clustering Complete LINK (HACCLINK) and Markovian scanning algorithm *i.e.*, Simulated ANNealing (SANN).

Intuition of the methodology: As introduced previously, we aim to construct a database of master templates. With this reduced database, an attacker has fewer templates to target but retains the ability to impersonate all individuals in the original database. To achieve this, we seek a list of templates that can represent all individuals. The idea is to group templates that are sufficiently close and calculate the master templates for these groups. Since the threshold is given by ε and the master template of the group must be able to impersonate all users within the group, we construct groups where the maximum distance between all members is at most 2ε . If the distance is greater than 2ε , it is evident that no point exists that is within ε of every other point in the group because there would be at least one point at least $\varepsilon + 1$ away from any potential master template. Another way of looking at it is that the intersection of all the ε -radius balls in the group would be empty. In summary, the methodology of this chapter can be described as follows: firstly, template groups are identified whose intersection balls are non-empty; secondly, a point is found in each intersection; and thirdly, these points are retained to constitute a new reduced database.

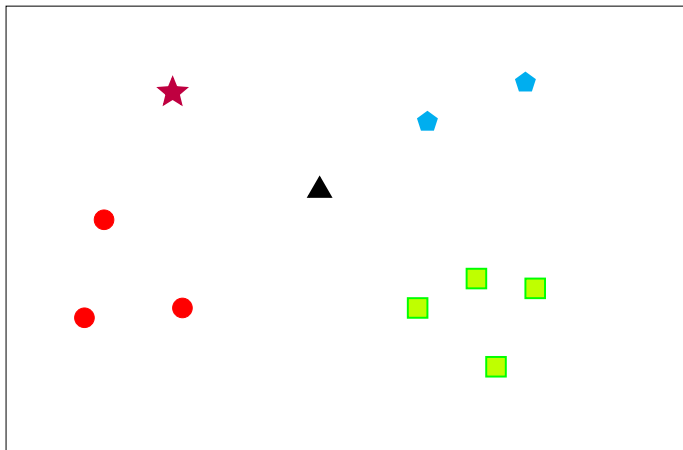


Figure 2.3 – Representation of a clustered database for a threshold of ϵ .

Agglomerative Clustering.

Consider $M_{\mathcal{B}}$ as the dissimilarity matrix of a template *database* \mathcal{B} , using the Hamming distance. $M_{\mathcal{B}}$ serves to compute template clusters, denoted by C_{ϵ} , where the distance between two templates within the same cluster is at most s . To achieve this clustering, an agglomerative clustering method is employed, which falls under the category of hierarchical clustering. This method involves iteratively merging the two closest groups of templates. Initially, there are $|\mathcal{B}|$ groups, each representing a single template, and the process continues until all groups are merged into one single cluster.

Standard post-processing is necessary to determine the appropriate iteration at which the algorithm should be terminated, ensuring the extraction of a relevant set of template clusters. However, we establish a termination condition such that the clustering algorithm stops when it becomes impossible to obtain template clusters that satisfy the required property outlined below:

$$\forall i \in \llbracket 1, n \rrbracket, \forall a, b \in C_i, d_H(a, b) \leq s.$$

The Agglomerative Clustering algorithm we used then corresponds to a slight variation of the HACCLINK (Hierarchical Agglomerative Clustering Complete LINK) presented by Defays [144].

By using the aforementioned clustering method, we obtain a set of template clusters, for which the inner-cluster distance suggests that there may exist at least one master-template for the cluster. An additional step is described below whose aim is to determine potential master-templates, if there exist some. Figure 2.3 depicts a clustered database where the distance between any two points within a given cluster is at most ϵ .

Master-Template of a Template Group.

The goal of this part is to compute the ϵ -master-template of each cluster as shown in Figure 2.4. We consider having a group of templates $L \subset D$ verifying

$$\forall i \in \llbracket 1, n \rrbracket, \forall a, b \in C_i, d_H(a, b) \leq s.$$

and for which we aim to find a master template. We emphasize that this problem can be formulated as a modified case of the Closest-String Problem (CSP) which is defined as follows.

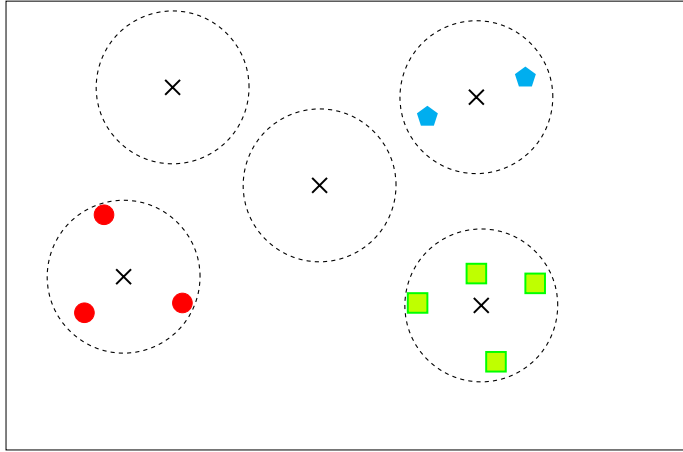


Figure 2.4 – The figure depicts the representation of a clustered database for a threshold of ε . The ε -master template of each cluster is represented by a cross.

Definition 2.2.3 (Modified closest-string problem (MCSP)). Given $S = \{s_1, s_2, \dots, s_m\}$ a set of strings with length n and d a distance, find a center string t of length m such that for every string s in S , $d_H(s, t) \leq d$.

The CSP is known as a *NP*-hard problem [139, 128, 129], and there exist algorithms to solve that kind of problem, see among others [125, 132, 99].

Definition 2.2.4 (Closest-String Problem (CSP)). Given $S = \{s_1, s_2, \dots, s_m\}$ a set of strings with length n , find a center string t of length m minimizing d such that for every string s in S , $d_H(s, t) \leq d$.

Given the relationship between the problems outlined in Definition 2.2.4 and 2.2.3, it can be asserted that the Modified Closest-String Problem (MCSP) is *NP*-hard.

Theorem 2.2.1 (MCSP is *NP*-hard). The Modified Closest-String Problem (MCSP) is *NP*-hard.

Proof. Let A be the oracle for the MCSP and (S) a CSP problem instance. Thus, on at most n calls to A , the closest-string problem can be solved. The solver B of the closest-string problem sends to A the following instances: $(S, 1), (S, 2), \dots, (S, i \leq n)$ and stops at the i -th instance for which A finally comes with the solution t . Then, B returns the pair (t, i) , a solution to the initial problem. Since B can be reduced in polynomial time to A and B is *NP*-hard, A is also *NP*-hard. The reduction is trivial in the other direction. ■

To our knowledge, the MCSP has not been previously tackled in the literature. As this problem is hard, employing brute force algorithms may not be efficient.

Let $\mathcal{B} = \{v_1, \dots, v_k\}$ represent a (TDB), and \mathcal{C} denote the ε -master-template set for D , which consists of ε -master-templates such that all points in D lie within a ball around a point in \mathcal{C} . The approach outlined below offers a constructive definition of the elements in \mathcal{C} , ensuring that $\mathcal{C} \neq \emptyset$. The following result emphasizes the relationship between \mathcal{C} and the balls $B_{\varepsilon, i} = \{u \in \mathbb{Z}_2^n \mid d_H(u, v_i) \leq \varepsilon\}$.

Proposition 2.2.1 (\mathcal{C} is the intersection of the balls of radius ε). Let $\mathcal{B} = \{v_1, \dots, v_k\}$ be a template database and \mathcal{C} the ε -master-template set for \mathcal{B} . Then, $\mathcal{C} = \bigcap_{i \in \{1, \dots, k\}} B_i$.

Proof. Let $p \in \bigcap_{u \in \mathcal{B}} B_\varepsilon(u)$. Then, $\forall u \in \mathcal{B}, p \in B_\varepsilon(u)$ with $B_\varepsilon(x) = \{y \in \mathbb{Z}_2^n \mid d_H(x, y) \leq \varepsilon\}$ the ball of radius ε center around x . Which implied that $\forall u \in \mathcal{B}, d_H(p, u) \leq \varepsilon$. And so, p is an ε -master-template for

\mathcal{B} . Then, $p \in \mathcal{C}$ which implies that $\cap_{u \in \mathcal{B}} B_\varepsilon(u) \subset \mathcal{C}$. Moreover, let $p \in \mathcal{C}$. Then, $\forall u \in \mathcal{B}, d_H(p, u) \leq \varepsilon$. So, $\forall u \in \mathcal{B}, p \in B_\varepsilon(u)$. Thus, $P \in \cap_{u \in \mathcal{B}} B_\varepsilon(u)$ and, $\mathcal{C} \subset \cap_{u \in \mathcal{B}} B_\varepsilon(u)$. Then, using both inclusion, $\mathcal{C} = \cap_{u \in \mathcal{B}} B_\varepsilon(u)$. ■

We denote by $p \in \mathcal{C}$ a master-template, and Proposition 2.2.1 indicates that determining all the master-template p is equivalent to determining the intersection of k -Hamming balls, which turn out to be formulated as the solutions of the following system:

$$d_H(p, v_i) \leq \varepsilon, \quad \forall i \in \{1, \dots, k\}. \quad (2.1)$$

With $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ two binary templates, the Hamming distance d_H between x and y can be rewritten as an inner product between $(n+2)$ -dimensional vectors in \mathbb{Z} , denoted X and Y . We have:

$$d_H(x, y) = \sum_{i=1}^n x_i + \sum_{i=1}^n y_i - 2 \sum_{i=1}^n x_i y_i = \langle X, Y \rangle$$

where

$$X = (x_1, x_2, \dots, x_n, 1, \sum_{i=1}^n x_i)$$

and

$$Y = (-2y_1, -2y_2, \dots, -2y_n, \sum_{i=1}^n y_i, 1).$$

Then, System 2.1 is a linear system, hence we can rely on a binary ILP (Integer Linear Programming) to solve it and then to compute \mathcal{C} . However, solving this system could be time-consuming in real-world cases since there are n parameters. Therefore, we suggest reducing System 2.1 by removing dependent variables and below are introduced necessary notations:

- For $K = \{k_1, \dots, k_{|K|}\} \subset \{1, \dots, n\}$, the Hamming distance over K is denoted by:

$$\forall u, v \in \mathbb{Z}_2^n, d_K = d_H((u_{k_1}, \dots, u_{k_{|K|}}), (v_{k_1}, \dots, v_{k_{|K|}})).$$

- Let $\mathcal{P}_D(K)$ a statement about $K \subset \{1, \dots, n\}$, $\mathcal{P}_B(K)$ holds if:

$$\forall u, v \in \mathcal{B}, d_K(u, v) \in \{0, |K|\}.$$

- The smallest partition $\{(K_1, \dots, K_{|I|}), K_i \subset \{1, \dots, n\} \mid \forall i \in \{1, \dots, |I|\}\}$ such that $\mathcal{P}_D(K_i)$ holds for all $i \in \{1, \dots, |I|\}$ is noted I . As I is the smallest possible partition, System 2.1 is reduced as much as it is possible.
- For $p \in \mathbb{Z}_2^n$ and $v \in \mathcal{B}$, $n_{v,i}$ denotes $d_{K_i}(p, v)$ and n_v^I denotes the parameters vector $(n_{v,1}, \dots, n_{v,|I|})$, written $N = (n_1, \dots, n_{|I|})$ for short when the context is clear.
- The distance vector $(d_H(v_1, v), \dots, d_H(v_{|B|}, v))$ is denoted by $d(v)$ with $v \in \mathcal{B}$ and $\mathcal{B} = (v_1, \dots, v_{|D|})$.

Then, with these notations, Theorem 2.2.2 can be stated, specifying a smaller version of System 2.1.

Theorem 2.2.2. For a given template database \mathcal{B} and for a given $v \in \mathcal{B}$, consider $L = \{p \in \mathbb{Z}_2^n \mid AN \leq \varepsilon - d(v)\}$ with $N = n_v^I$, $\varepsilon = (\varepsilon, \dots, \varepsilon)^T$, $n_{v,i}$ denotes $d_{K_i}(p, v)$, n_v^I denotes the parameters vector $(n_{v,1}, \dots, n_{v,|I|})$

and $A = (a_{i,j})$ a matrix of size $|I| \times |\mathcal{B}|$ whose the $(i, j)^{th}$ element is

$$a_{i,j} = \begin{cases} 1 & \text{if } d_{K_j}(v_1, v_i) = 0 \\ -1 & \text{if } d_{K_j}(v_1, v_i) = |K_j| \end{cases}$$

Then, $L = \mathcal{C}$ the ε -master-template set for \mathcal{B} .

Proof. Let \mathcal{B} be a database, $u \in \mathcal{B}$ a template, $p \in \mathbb{Z}_2^n$ a template and, $A_K(u) = \{v \in \mathcal{B} | d_K(u, v) = 0\}$. There are two cases:

1. If $v \in A_K(u)$ then, $d_K(u, v) = 0$.
2. Else, $v \in A_K(u)^c$ then, $d_K(u, v) = |K|$.

If $v \in A_K(u)^c$ then, $d_K(v, p) = |K| - d_K(p, u)$. However, as I is a partition of $\{1, \dots, n\}$, $d_H(u, p) = \sum_{K \in I} d_K(u, p)$.

Suppose that $p \in \mathcal{C}$ the ε -master-template set for D . As $\max_{u \in D} d_H(u, p) \leq \varepsilon$ then, $\sum_{K \in I} d_K(u, p) \leq \varepsilon$. Thus, for $v \in \mathcal{B}$, $d_K(v, p) = d_K(p, u)1_{A_K(u)}(v) + (|K| - d_K(u, p))1_{A_K(u)^c}(v)$. Then, for a given couple (u, v) , we have:

$$\begin{aligned} \sum_{K \in I} d(v, p) &= \sum_{K \in I} d_K(p, u)1_{A_K(u)}(v) + (|K| - d_K(u, p))1_{A_K(u)^c}(v) \\ &= \left(\sum_{K \in I} d_K(p, u) (1_{A_K(u)}(v) - 1_{A_K(u)^c}(v)) \right) + \sum_{K \in I} |K| 1_{A_K(u)^c}(v) \end{aligned}$$

Moreover, $d_H(u, v) = \sum_{K \in I} |K| 1_{A_K(u)^c}(v)$ then,

$$\sum_{K \in I} d(v, p) = \sum_{K \in I} d_K(p, u) (1_{A_K(u)}(v) - 1_{A_K(u)^c}(v)) + d_H(u, v).$$

Then,

$$\begin{aligned} &\sum_{K \in I} d_K(v, p) \leq \varepsilon \\ \Leftrightarrow &\sum_{K \in I} d_K(p, u) (1_{A_K(u)}(v) - 1_{A_K(u)^c}(v)) \leq \varepsilon - d_H(u, v) \\ \Leftrightarrow &A(u) d_K(p, u) \leq \varepsilon - d_H(u, v) \\ \Leftrightarrow &p \in L \end{aligned}$$

■

As I is required to reduce System 2.1, we assure with Lemma 2.2.1 that $I \neq \emptyset$, whatever the configuration of the set \mathcal{B} is.

Lemma 2.2.1 (I is not empty). $\forall \mathcal{B} \subset \mathbb{Z}_2^n$ such that $|\mathcal{B}| > 1$, $I \neq \emptyset$.

Proof. Let $\mathcal{B} \subset \mathbb{Z}_2^n$ be a template database such that $|\mathcal{B}| > 1$ and $K_i = \{i\}, \forall i \in \{1, \dots, n\}$. Therefore, $\sqcup_{i \in \{1, \dots, n\}} K_i = \{1, \dots, n\}$ and, $\forall i \in \{1, \dots, n\}$, $\mathcal{P}_{\mathcal{B}}(K_i) = \text{True}$. Then, $I = \{K_i, \forall i \in \{1, \dots, n\}\} \neq \emptyset$. ■

In the same vein, we have $|I| \leq n$. As $|I|$ corresponds to the number of parameters, the system described in Theorem 2.2.2 is always smaller or equivalent to System 2.1.

Example 2.2.2.1. Let $D = \{(1, 0, 1, 1, 0, 1, 1, 0), (1, 0, 0, 1, 0, 1, 1, 1), (1, 0, 1, 1, 1, 1, 1, 0), (1, 0, 0, 1, 1, 1, 0, 1)\}$ be a database represented as a matrix with the templates in rows. The identical or opposite columns are labeled with the same symbol, as follows:

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------|---|---|---|---|---|---|---|---|
| v_1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| v_2 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |
| v_3 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 0 |
| v_4 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 1 |
| | ♠ | ♠ | ■ | ♠ | ★ | ♠ | ▲ | ■ |

We remind that $\mathcal{P}_B(K)$ holds if, for all templates of B , their pairwise distance is $|K|$ or 0. Let $K = \{1, 2, 4, 6\}$ be the set of columns marked with a ♠. Then, $\mathcal{P}_B(K)$ holds. However, for $K = \{3, 7\}$, $\mathcal{P}_B(K)$ does not hold. If K is uniquely comprised of columns having the same symbol, the statement holds. If the columns that are identical or opposite are merged, the property $\mathcal{P}_B(K)$ holds. Finally, in this example, the partition I is $\{\underbrace{\{1, 2, 4, 6\}}_{♠}, \underbrace{\{3, 8\}}_{■}, \underbrace{\{5\}}_{★}, \underbrace{\{7\}}_{▲}\}$.

Theorem 2.2.2 suggests that determining the ε -master-template set for B , which involves finding the intersection of $|\mathcal{B}|$ balls in \mathbb{Z}_n^2 , can be simplified to solving a potentially small linear system. Although this system can be solved efficiently with powerful tools such as GUROBI [95], we advocate for the use of simpler algorithms in this scenario. Depending on the configuration of B , it may be feasible to derive a linear system such that it becomes straightforward to determine the space of potential solutions and find a solution using any Markovian scanning algorithm. If \mathcal{N} denotes the set of the possible solutions N for the linear system described in Theorem 2.2.2, we have:

$$\mathcal{N} = \prod_{k=1}^{|\mathcal{I}|} \{0, \dots, \min(\varepsilon, |K_k|)\}$$

since, for $k \in \{1, \dots, |\mathcal{I}|\}$, $n_{v,k}$ corresponds to the distance $d_{K_k}(v_k, v)$, which can not be greater than $|K_k|$, and in the other hand if $d_{K_k}(v_k, v) > \varepsilon$ then, N does not belong to L . Note that depending on the dimension of \mathcal{N} , finding a solution N can be efficiently accomplished either through a brute force algorithm, especially for small-dimensional sets \mathcal{N} , or via a more parsimonious algorithm when dealing with high-dimensional sets. As the dimension of \mathcal{N} depends among other factors on B , we consider that the use of one of both approaches should be determined with regard to practical context-specific consideration. In the following, we assume that the solution set \mathcal{N} is a high dimensional set. To tackle this problem in high dimension, we propose to rely on an efficient and simple algorithm: the Simulating Annealing algorithm [142]. Below we detail features of the Simulating Annealing algorithm that we tune to obtain good performances in our numerical evaluations. It is composed of the following parameters:

- *Energy*: We define the following energy in such a way that as it increases, the closer N gets to solving the system: $E(N) = \sum_{i=1}^{|\mathcal{I}|} f((\varepsilon - d(v) - AN)_i)$ where f is a ReLU type function: $f(x) = \min(0, x)$.
- *Cooling Schedule*: In practice, we observe that the solution-finding process is not significantly affected by the cooling function, as shown in Table 2.1. Therefore, we suggest using a linearly decreasing temperature function. The initial temperature is set such that in the very first iterations, all potential moves are accepted regardless of the chosen initial point.

| Cooling method | n | ε | Number of clients | Error in % | Time (ms) |
|-----------------------|-----|---------------|-------------------|------------|-----------|
| Additive cooling | 45 | 10 | 50 | 0.1 | 7 |
| | 50 | 10 | 50 | 0.6 | 11 |
| | 55 | 10 | 50 | 3.1 | 9 |
| | 60 | 10 | 50 | 8.3 | 11 |
| | 65 | 10 | 50 | 47.2 | 19 |
| Linear multiplicative | 45 | 10 | 50 | 0.6 | 12 |
| | 50 | 10 | 50 | 0.8 | 11 |
| | 55 | 10 | 50 | 3.7 | 5 |
| | 60 | 10 | 50 | 5.3 | 16 |
| | 65 | 10 | 50 | 35.3 | 18 |
| Exponential | 45 | 10 | 50 | 0.2 | 7 |
| | 50 | 10 | 50 | 1.2 | 10 |
| | 55 | 10 | 50 | 4.3 | 3 |
| | 60 | 10 | 50 | 6.9 | 4 |
| | 65 | 10 | 50 | 40.8 | 15 |
| Logarithmic | 45 | 10 | 50 | 0.5 | 6 |
| | 50 | 10 | 50 | 1.4 | 3 |
| | 55 | 10 | 50 | 2.9 | 3 |
| | 60 | 10 | 50 | 6.6 | 10 |
| | 65 | 10 | 50 | 40.8 | 10 |

Table 2.1 – Comparison of cooling methods for our simulated annealing.

- *Proposal distribution*: According to computational considerations and for the sake of numerical performance, we define a proposal distribution for which the support is the neighbor set. Moreover, we choose a non-symmetric proposal that preferentially promotes neighbors that increase energy.
- *Termination*: The algorithm is terminated either it reaches the maximum iteration number (about 200k iterations), or if a solution is found, which corresponds to a vector N with null energy.

Database Partitioning Algorithm.

With the aforementioned building blocks, we propose Algorithm 1 to partition the template database. It takes as inputs \mathcal{B} a template database and a threshold ε and returns an ε -master-template set denoted by MTS.

Figure 2.5 provides a schematic illustration of the database partitioning algorithm to provide an intuitive understanding of its underlying mechanism.

Evaluation of the Database Partitioning Algorithm

To model the worst case for an attacker, the templates are randomly drawn from \mathbb{Z}_2^n . The parameters selected for analysis were specifically chosen to illustrate a change in the number of clusters and computation time, allowing for the observation of the algorithm's behavior. For each configuration, experimentations are replicated 1000 times, averaged and presented in Table 2.2 with the following notations:

- n : the space dimension,

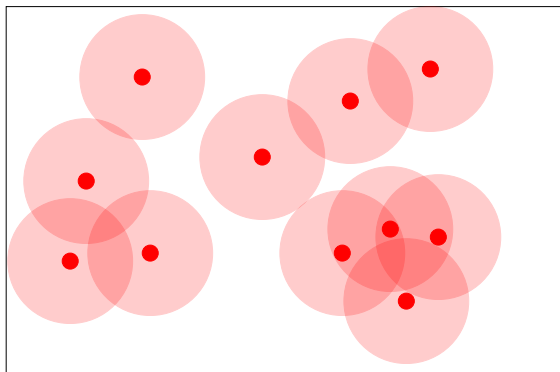
Algorithm 1 : Database partitioning algorithm

Input : \mathcal{B}, ϵ
Output : MTS

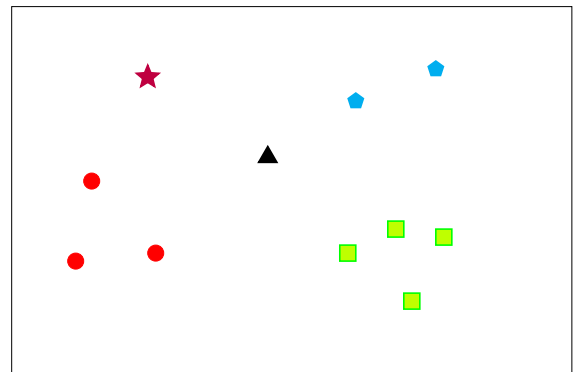
```

1 Function Partitioning( $y$ ):
2   Set  $s$  to  $2\epsilon$ .
3   Set MTS to [ ].
4   while  $D \neq \emptyset$  do
5     Compute cluster  $C$ s using  $D$  and  $s$ .
6     foreach cluster  $c$  in  $C$ s do
7       Search the master template  $t$  for  $c$ .
8       if a cover template  $t$  is found for  $c \in C$  then
9         Set  $D$  to  $D \setminus c$  and add  $t$  to MTS.
10      end
11      Set  $s$  to  $s - 1$ .
12    end
13  end
14 return MTS

```



(a) Initial database.



(b) Clustering.

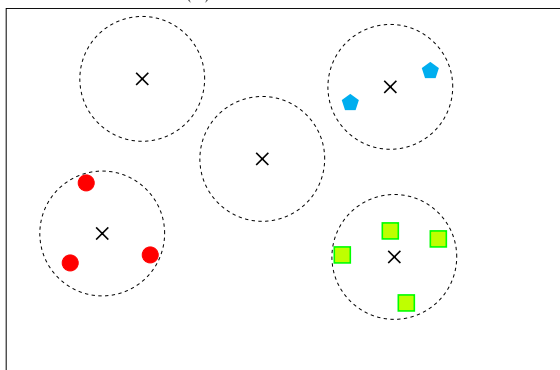
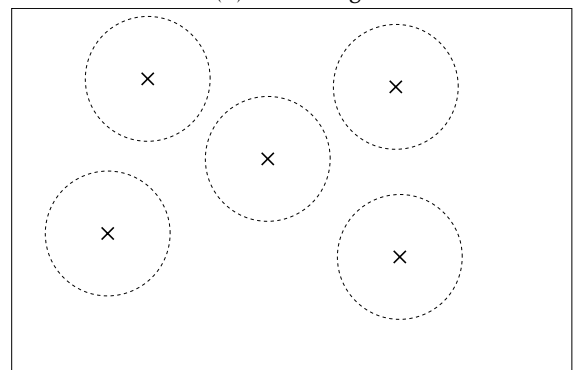
(c) Find the ϵ -master template of each cluster.(d) Remaining database or ϵ -master template set.

Figure 2.5 – The figure depicts the representation of the database partitioning algorithm.

| n | ϵ | #clients | # clust | Time (s) |
|-----|------------|----------|---------|----------|
| 15 | 10 | 50 | 2.000 | 0.310 |
| 20 | 10 | 50 | 8.000 | 0.345 |
| 25 | 10 | 50 | 12.659 | 34.588 |
| 30 | 10 | 50 | 29.000 | 70.683 |
| 35 | 10 | 50 | 44.000 | 0.629 |
| 40 | 10 | 50 | 45.000 | 0.514 |
| 45 | 10 | 50 | 49.000 | 0.488 |
| 35 | 7 | 50 | 50.0 | 0.343 |
| 35 | 9 | 50 | 48.0 | 0.687 |
| 35 | 11 | 50 | 41.0 | 1.017 |
| 35 | 13 | 50 | 28.0 | 1.142 |
| 35 | 15 | 50 | 9.0 | 0.933 |
| 35 | 10 | 30 | 28.000 | 0.090 |
| 35 | 10 | 50 | 44.000 | 0.193 |
| 35 | 10 | 70 | 57.000 | 0.336 |
| 35 | 10 | 90 | 69.000 | 0.485 |
| 35 | 10 | 110 | 82.000 | 0.700 |
| 35 | 10 | 130 | 95.000 | 5.140 |
| 35 | 10 | 150 | 112.000 | 5.796 |
| 35 | 10 | 170 | 119.000 | 9.753 |

Table 2.2 – Summary of the experiments of the space partitioning algorithm.

- ϵ : the threshold,
- $\#clients$: the number of templates in the TDB,
- $\#clust$: the number of clusters found with Algorithm 1,
- $Time$ is the running time of Algorithm 1,

When n is increasing, the template space expands. We observe that as the space expands, the computation time increases with the number of clusters. This is because the templates tend to move farther apart in this case. The more clusters there are, the longer the computation time. However, there is a point where the computation time drastically reduces. This happens because as the templates spread out, clusters containing only one template become predominant, and finding the master template of these clusters is trivial. When the number of clients $\#clients$ increases, the density of the space increase. In this case, the computation time increases. This is easily explained by the fact that more clients mean more time is required to compute the pairwise distance matrix used for the clustering. Additionally, we observe an increase in the number of clusters. Given the increased density of the space with more clients, there is no trivial master template because fewer clusters may contain only one template. When ϵ increases, the size of the balls grows. As the size of the balls increases, it is natural to observe a decrease in the number of clusters and an increase in computational time, as finding a center for each cluster becomes necessary. However, as the radius size continues to grow, the computation time should decrease again because there are fewer centers to find, and it becomes simpler to find them since more points fall within the intersection of the templates in the clusters.

Prevent this Attack. The most effective approach to mitigating this attack is to minimize the occurrence of near-collisions. This is achieved by adjusting the database size according to the

results presented in Section 2.2.3.

2.2.3 Multi-Near-Collisions and Master-templates

This section introduces the concept of near-collisions, which is examined in further detail in the subsequent sections. In essence, a near-collision occurs when two individuals in the same database possess attributes that allow them to authenticate as one another. Bounds on near-collisions are defined in function of template size, decision threshold, and database population. We then revisit the master templates, as discussed in the previous section, in order to highlight the probability of their occurrences. Finally, we conduct numerical evaluations on near-collisions, with an emphasis on the study of critical populations. The critical population is defined as the number of individuals at which the probability of a near-collision exceeds a given threshold.

Intuition of the Population Size Impact on Near-Collision: As the population size increases, the likelihood of near-collisions between templates also increases. This is because a larger population results in a higher density of templates within the metric space. Consequently, the distance between some templates inevitably decreases, leading to a greater probability of two templates being within each other's ε -radius. There exists a critical population size where the probability of near-collisions exceeds a fixed threshold. At this point, the risk of impersonation becomes significant, as two different users might have templates close enough to be considered matches.

Definitions

The database \mathcal{B} contains N templates which are vectors distributed in \mathbb{Z}_2^n . This distribution is considered uniform if templates result from a salted or a secret-based transformation. If the sequence of treatments applied to the feature vectors is deterministic, the templates can be regarded as non-uniformly distributed in the set of integers modulo two of any given dimension. The following definitions are required:

Definition 2.2.5 (Strong ε -near-collision for t). *For a secret template t of \mathcal{B} , a strong near-collision occurs if there is another template $a \in \mathcal{B}$ such that $d_{\mathcal{H}}(t, a) \leq \varepsilon$.*

Notice that for a given secret template t of a database \mathcal{B} , the probability of $s\text{-nc}_t(\varepsilon, N)$ the event "At least one of the $N - 1$ other users of \mathcal{B} matches the given template t " is $\mathbb{P}(s\text{-nc}_t(\varepsilon, N)) = 1 - (1 - V_\varepsilon)^{N-1}$. In other words, in the case of this definition, the occurrence of a strong ε -near-collision geometrically increases, if a given targeted template t is considered. This scenario corresponds to a simple case for which we can easily obtain an interpretable result. However, it is only an intermediate step, since this scenario does not accurately represent the case of a realistic near-collision. In the following, we focus on a more general case by not considering a given targeted template, which involves different near-collision events (Definitions 2.2.6 and 2.2.7) for which it is not possible to smoothly derive probability and complexity results.

Definition 2.2.6 (Weak ε -near-collision). *For \mathcal{B} a biometric database, a weak ε -near-collision is occurring in \mathcal{B} if there exists two templates $a, b \in \mathcal{B}$ such that $d_{\mathcal{H}}(a, b) \leq \varepsilon$.*

In other words, a weak near-collision occurs if there exists a pair of templates a and b in the secret database such that a (resp. b) is inside the ball of center b (resp. a). A representation of a near-collision is given in Figure 2.6.

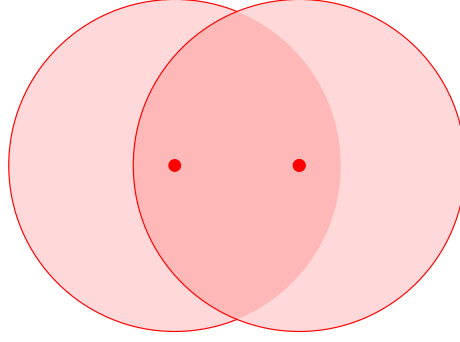


Figure 2.6 – Illustration of near-collisions of two templates.

Note that if a and b weak near-collides, then we have two strong near-collisions, for a , and for b . This definition can be generalized to the case of multi-near-collisions.

Definition 2.2.7 (Weak (m, ε) -near-collision). For \mathcal{B} a biometric database, an (m, ε) -near-collision with $m \geq 2$ is occurring if there exists m templates $a_1, \dots, a_m \in \mathcal{B}$ such that for all $i, j \in \{1, \dots, m\}$, $d_{\mathcal{H}}(a_i, a_j) \leq \varepsilon$.

Risk Related to Weak Collisions

For what follows, W_N denotes the event that a weak collision is found in a set of N enrolled templates. A particular distinction concerning this case is that we have to compare the matches between all the pairs in the database. The probability calculus then does not follow directly as before but comes down to a problem similar to the Birthday Problem. Since W_N is an event whose probability is difficult to calculate, the theorem 2.2.3 below provides upper and lower bounds for this probability.

Theorem 2.2.3. For a database \mathcal{B} of uniformly drawn templates, the probability of W that there is a weak collision for at least two templates is bounded as follows

$$1 - (1 - V_\varepsilon)^{N(N-1)/2} \leq \mathbb{P}(W) \leq 1 - \prod_{j=1}^N \max(0, 1 - (j-1)V_\varepsilon)$$

where V_ε is the measure of an ε -ball.

Proof. Consider $\mathcal{B} = (v_1, \dots, v_n)$ the database. Notice that

$$\mathbb{P}(W) = \mathbb{P}\left(\exists u, v \in \mathcal{B}, u \in B_\varepsilon(v)\right) = 1 - \mathbb{P}\left(\forall u, v \in \mathcal{B}, u \notin B_\varepsilon(v)\right),$$

and the term can be developed as:

$$\mathbb{P}\left(\forall u, v \in \mathcal{B}, u \notin B_\varepsilon(v)\right) = \prod_{j=2}^N \mathbb{P}\left(v_j \notin \bigcup_{k=1}^{j-1} B_\varepsilon(v_k) \mid \overline{W_{j-1}}\right). \quad (2.2)$$

Notice that if the $j-1$ first ε -balls are disjoint, then the cardinal of the complementary of $\bigcup_{k=1}^{j-1} B_\varepsilon(v_k)$ is minimal. Then, by denoting D_{j-1} the event that the $j-1$ first ε -balls are disjoint, a lower bound of the conditional probability is:

$$\mathbb{P}\left(v_j \notin \bigcup_{k=1}^{j-1} B_\varepsilon(v_k) \mid \overline{W_{j-1}}\right) \geq \mathbb{P}\left(v_j \notin \bigcup_{k=1}^{j-1} B_\varepsilon(v_k) \mid D_N\right) = 1 - (j-1)V_\varepsilon.$$

Furthermore, remark that if j is too large, then D_{j-1} cannot occur and in the same vein for $v_j \notin \bigcup_{k=1}^{j-1} B_\varepsilon(v_k)$. As a consequence, in this case, the targeted conditional probability is necessarily equal to 0. Merging both cases leads to the upper bound of $\mathbb{P}(W)$.

For the lower bound, Equation (2.2) can be rewritten as follows by considering that the templates are not in weak near-collisions is equivalent to considering that each pair of templates is not in weak near-collision:

$$\mathbb{P}\left(\forall u, v \in \mathcal{B}, u \notin B_\varepsilon(v)\right) = \mathbb{P}\left(v_i \notin B_\varepsilon(v_j), \text{ for } i \neq j \in \{1, \dots, N\}\right)$$

For what follows, we consider a path in the set of pairs c , those the k^{th} element is denoted by c_k and it consists of the values c_k^1 and c_k^2 . By using the chain rule:

$$\mathbb{P}\left(v_i \notin B_\varepsilon(v_j), \text{ for } i \neq j \in \{1, \dots, N\}\right) = \mathbb{P}\left(v_{c_1^2} \notin B_\varepsilon(v_{c_1^1})\right) \prod_{k=2}^{|c|} \mathbb{P}\left(v_{c_k^2} \notin B_\varepsilon(v_{c_k^1}) \mid v_{c_j^2} \notin B_\varepsilon(v_{c_j^1}), \forall j < k\right)$$

where $|c| = \binom{N}{2}$. First, recall that $\mathbb{P}\left(v_{c_1^2} \notin B_\varepsilon(v_{c_1^1})\right) = 1 - V_\varepsilon$. Next, remark that in the conditional part of the conditional probability, if for all j , $v_{c_j^1}$ and $v_{c_j^2} \notin \{v_{c_k^1}, v_{c_k^2}\}$ then

$$\mathbb{P}\left(v_{c_k^2} \notin B_\varepsilon(v_{c_k^1}) \mid v_{c_j^2} \notin B_\varepsilon(v_{c_j^1}), \forall j < k\right) = \mathbb{P}\left(v_{c_k^2} \notin B_\varepsilon(v_{c_k^1})\right) = 1 - V_\varepsilon.$$

Otherwise, for some events, there is a least one value j with $v_{c_j^1}$ or $v_{c_j^2} \in \{v_{c_k^1}, v_{c_k^2}\}$. Then the conditional event imposes a constraint on the possible templates in \mathbb{Z}_2^n for $v_{c_k^1}$ and/or $v_{c_k^2}$, and we denote by $R_{k,1}$ and $R_{k,2}$ the templates that cannot be taken respectively by $v_{c_k^1}$ and $v_{c_k^2}$. The event $v_{c_k^2} \notin B_\varepsilon(v_{c_k^1})$ can be modeled by fixing a template $v_{c_k^1} \in \mathbb{Z}_2^n \setminus R_{k,1}$ and next by choosing $v_{c_k^2}$ among $(\mathbb{Z}_2^n \setminus R_{k,2}) \cap \overline{B_\varepsilon(v_{c_k^1})}$, thus

$$\begin{aligned} \mathbb{P}\left(v_{c_k^2} \notin B_\varepsilon(v_{c_k^1}) \mid v_{c_j^2} \notin B_\varepsilon(v_{c_j^1}), \forall j < k\right) &= \mathbb{P}\left(v_{c_k^2} \in (\mathbb{Z}_2^n \setminus R_{k,2}) \cap \overline{B_\varepsilon(v_{c_k^1})} \mid v_{c_k^1} \in \mathbb{Z}_2^n \setminus R_{k,1}\right) \\ &\leq \mathbb{P}\left(v_{c_k^2} \notin B_\varepsilon(v_{c_k^1}) \mid v_{c_k^1} \in \mathbb{Z}_2^n \setminus R_{k,1}\right) = 1 - V_\varepsilon. \end{aligned}$$

It follows that the probability of a least one weak near-collision can be bounded below by $1 - (1 - V_\varepsilon)^{N(N-1)/2}$. ■

To highlight the significance of this result, we present Proposition 2.2.2 below, demonstrating the asymptotic proximity of these bounds to each other. To ease the equations, b_{inf} and b_{sup} respectively denote the lower and upper bounds of $\mathbb{P}(\overline{W})$, according to Theorem 2.2.3. In particular, notice that asymptotic results about the gap between b_{inf} and b_{sup} derive from different cases of ratio between N and V_ε , where V_ε is related to n as $V_\varepsilon = |B_\varepsilon|/2^n$.

Proposition 2.2.2. *For N the number of users in the database and n the template size. Let α be a parameter to define N as a portion of the number of possible templates: $NV_\varepsilon = \alpha$.*

- (a) *If $\log_2 N \leq \frac{c}{2}n$, with $0 \leq c < 1$ then $b_{\text{inf}} \sim b_{\text{sup}}$ as $n \rightarrow +\infty$.*
- (b) *If $\alpha < 1$ then, $\exists 0 < k \leq 1$ so that $kb_{\text{sup}}^2 \leq b_{\text{inf}} \leq b_{\text{sup}}$ for all N .*
- (c) *If $\alpha > 1$ then $b_{\text{inf}} = 0$, $2^n = O(N)$ and $b_{\text{sup}} = o(2^{-n})$ as $n \rightarrow +\infty$.*

Proof. First, rewrite b_{inf} and b_{sup} and provide bounds:

$$b_{\text{sup}} = (1 - V_\varepsilon)^{N(N-1)/2} = \exp\left(\frac{N(N-1)}{2} \log(1 - V_\varepsilon)\right) \leq \exp\left(-NV_\varepsilon \frac{(N-1)}{2}\right) \quad (2.3)$$

Next note that $b_{\text{inf}} = 0$ for $NV_\varepsilon \geq 1$, so the lower bound for b_{inf} is only given for $NV_\varepsilon < 1$, and in this case $\alpha = NV_\varepsilon$. For what follows, we recall that the gamma function can be written as [101]:

$$\forall x, \exists 0 < \theta < 1, \Gamma(x+1) = \sqrt{2\pi} x^{x+1/2} \exp\left(-x + \frac{\theta}{12x}\right), \quad (2.4)$$

and then

$$\begin{aligned} b_{\text{inf}} &= \prod_{j=1}^N \max(0, 1 - (j-1)V_\varepsilon) = \prod_{j=1}^N 1 - (j-1)V_\varepsilon = V_\varepsilon^N \frac{\Gamma(V_\varepsilon^{-1} + 1)}{\Gamma(V_\varepsilon^{-1}(1-\alpha) + 1)} \\ &= \exp\left(N \log V_\varepsilon - V_\varepsilon^{-1} + (V_\varepsilon^{-1} + \frac{1}{2}) \log V_\varepsilon^{-1} + \frac{\theta_1}{12V_\varepsilon^{-1}} + \right. \\ &\quad \left. V_\varepsilon^{-1}(1-\alpha) - (V_\varepsilon^{-1}(1-\alpha) + \frac{1}{2}) \log(V_\varepsilon^{-1}(1-\alpha)) - \frac{\theta_2}{12V_\varepsilon^{-1}(1-\alpha)}\right) \\ &= \exp\left(-N \left(1 + \frac{1-\alpha}{\alpha} \log(1-\alpha)\right) - \frac{1}{2} \log(1-\alpha) + \frac{\alpha\theta_1}{12N} - \frac{\alpha\theta_2}{12N(1-\alpha)}\right) \\ &\geq \exp\left(-N \left(1 + \frac{1-\alpha}{\alpha} \log(1-\alpha)\right) - \frac{1}{2} \log(1-\alpha) - \frac{\alpha}{12N(1-\alpha)}\right) \end{aligned} \quad (2.5)$$

Case (a): As $\log_2 N \leq \frac{\varepsilon}{2}n$, then $\alpha = NV_\varepsilon < 1$ and then $\max(0, 1 - (j-1)V_\varepsilon) \neq 0$ for all $j \in \{1, \dots, N\}$. Then, with (2.3) and (2.5), we obtain the upper bound:

$$\frac{b_{\text{sup}}}{b_{\text{inf}}} < \exp\left(N \left(1 - \frac{\alpha}{2} + \frac{1-\alpha}{\alpha} \log(1-\alpha)\right) + \frac{1}{2} (\alpha + \log(1-\alpha)) + \frac{\alpha}{12N(1-\alpha)}\right).$$

Since $\log(1-\alpha) \leq -\alpha$, then $\left(1 - \frac{\alpha}{2} + \frac{1-\alpha}{\alpha} \log(1-\alpha)\right) \leq \alpha/2$ and $\alpha + \log(1-\alpha) \leq 0$, so that

$$\frac{b_{\text{sup}}}{b_{\text{inf}}} < \exp\left(\frac{\alpha}{2}N + \frac{\alpha}{12N(1-\alpha)}\right).$$

With $\log_2 N \leq \frac{\varepsilon}{2}n$, we have that $\frac{\alpha}{2}N = \frac{1}{2}N^2V_\varepsilon = N^2|B_\varepsilon|2^{-n-1} \leq |B_\varepsilon|2^{-n(1-c)-1}$, which ensures that $\frac{\alpha}{2}N$ converges to 0 as n increases, since $c < 1$. It also ensures that $\frac{\alpha}{12N(1-\alpha)}$ tends to 0, and then the result follows.

Case (b): Recall that $NV_\varepsilon = \alpha < 1$. To ease the following equations, $r_\alpha = 1 + \frac{1-\alpha}{\alpha} \log(1-\alpha)$ and $s_\alpha = -\frac{1}{2} \log(1-\alpha) - \frac{\alpha}{12N(1-\alpha)}$. Similar to Case (a) remark that, for all $\alpha < 1$, $r_\alpha \leq \alpha$, then,

$$b_{\text{inf}} \geq e^{-N\alpha+s_\alpha} \geq e^{s_\alpha-\alpha} b_{\text{sup}}^2$$

It remains that s_α depends on N and in order to obtain k , consider the following bound for s_α :

$$s_\alpha \geq -\frac{1}{2} \log(1-\alpha) - \frac{1}{12} \times \frac{\alpha}{1-\alpha}.$$

Case (c): If $\alpha > 1$ and n increases to $+\infty$ then N also increases to $+\infty$, ince $N > 2^n/|B_\varepsilon|$. It follows

that $\exp\left(-\alpha \frac{(N-1)}{2}\right) \rightarrow 0$, which ensures that $b_{\text{sup}} \rightarrow 0$. ■

This result shows how close these bounds are when n becomes large. These bounds help interpret the risk of the system due to the presence of a weak collision. Specifically, the lower bound (which corresponds to the upper bound of Theorem 2.2.3) is particularly interesting because it represents a pessimistic view of the security of the system concerning this issue. Furthermore, the interpretation of this asymptotic result is facilitated by the two results presented below. Proposition 2.2.3 adapts a classical result from the Birthday Problem to the context of a weak collision, providing the number of users required for the probability of a weak collision to be at least one-half. Then, Proposition 2.2.4 shows how to calibrate the parameter ε (through the value of V_ε) to ensure that the probability of a weak collision is below a risk threshold r .

Proposition 2.2.3. *In the same context as Theorem 2.2.3, $\mathbb{P}(W) \geq r$ with $0 < r < 1$ if*

$$N_r \geq \frac{1}{2} + \sqrt{\frac{1}{4} + 2 \frac{\log(1-r)}{\log(1-V_\varepsilon)}}.$$

Proof. According to Theorem 2.2.3, solving $\mathbb{P}(W) \geq r$ is reduced to solving $1 - (1 - V_\varepsilon)^{N_r(N_r-1)/2} \geq r$, which is similar to solve

$$N_r^2 - N_r - 2 \frac{\log(1-r)}{\log(1-V_\varepsilon)} \geq 0$$

and determining the only positive solution for N_r leads to the result. ■

The classical application of the birthday problem takes $r = \frac{1}{2}$ yielding

$$N_{\frac{1}{2}} \geq \frac{1}{2} + \sqrt{\frac{1}{4} + 2 \frac{\log \frac{1}{2}}{\log(1-V_\varepsilon)}}.$$

For the next result, we assume $NV_\varepsilon \leq 1 - \gamma$, where γ can be chosen close to 0. This assumption greatly simplifies the equations and improves the readability of the result. This assumption is reasonable because if V_ε becomes close to 1, then the probability of a weak collision is high.

Proposition 2.2.4. *Let $r > 0$ a risk level and assume that $NV_\varepsilon \leq 1 - \gamma$, then $\mathbb{P}(W) \leq r$ if*

$$\log V_\varepsilon \leq -\log\left(N(N-1) + \frac{1}{12\gamma}\right) + \log(-\log(1-r)).$$

Proof. First, with regards to Theorem 2.2.3, notice that $\mathbb{P}(W) \leq r$ if $\prod_{j=1}^N 1 - (j-1)V_\varepsilon \geq 1 - r$, since $NV_\varepsilon \leq 1 - \gamma$ implies that $\max(0, 1 - (j-1)V_\varepsilon) = 1 - (j-1)V_\varepsilon$. By following the same procedure as in proof of Proposition 2.2.2, this reduces to solve

$$-N\left(1 + \frac{1 - NV_\varepsilon}{NV_\varepsilon} \log(1 - NV_\varepsilon)\right) - \frac{1}{2} \log(1 - NV_\varepsilon) - \frac{V_\varepsilon}{12(1 - NV_\varepsilon)} \geq \log(1 - r).$$

Since $\log(1 - NV_\varepsilon) \leq -NV_\varepsilon$ and $NV_\varepsilon \leq 1 - \gamma$, this leads to

$$-N^2V_\varepsilon + \frac{1}{2}NV_\varepsilon - \frac{V_\varepsilon}{12\gamma} \geq \log(1 - r),$$

and then the result is obtained by isolating V_ε . ■

As an illustration of Proposition 2.2.4, consider a system with $N = 10^{10}$ users for which the template size is $n = 128$. We assume that $\gamma = 0.001$, which is low. Notice that $\gamma = 0.001$ holds if $\varepsilon \leq 28$, and if $\varepsilon = 29$ then $\mathbb{P}(W) = 1$ according to the Pigeonhole Principle. Then, for $r = 0.001$ with regards to Proposition 2.2.4, $\mathbb{P}(W) \leq 0.001$ if $\varepsilon \leq 11$. In the same setting for which we target to have $\mathbb{P}(W) \leq 0.001$, if $N = 10^6$ then it requires $\varepsilon \leq 20$ and if $N = 10^3$ it requires $\varepsilon \leq 30$.

2.2.4 Risk related to Master-templates

The previous section focuses on the probability of a weak collision, which poses a potential security flaw. In this section, we examine master templates, which are defined below and specify possible weak collision cases to assess the security risks.

Definition 2.2.8 ((k, ε) -master template). *A (k, ε) -master template t with respect to \mathcal{B} is a template such that there exists a k distinct templates a_1, \dots, a_k in \mathcal{B} with $d_{\mathcal{H}}(a_i, t) \leq \varepsilon$ for all $i \in \{1, \dots, k\}$.*

To compute the probability that there exists a (N, ε) -master template, we firstly introduce $\mathcal{C}(\varepsilon, N)$ the number of different sets of N vectors that can be strictly inside a given ε -ball,

$$\mathcal{C}(\varepsilon, N) = \left\{ \mathcal{B} = (v_1, \dots, v_N) \in (\mathbb{Z}_2^n)^N, \mathcal{B} \subset B_\varepsilon \right\},$$

and $\mathcal{C}_v(\varepsilon, N)$ is the same number but for a ε -ball centered on a given template v . Then, the following theorem provides the probability that there exists a (N, ε) -master template.

Theorem 2.2.4. *For a template database \mathcal{B} of size N , the probability that there exists a (N, ε) -master template is:*

$$\mathbb{P}(\mathcal{B} \subset B_\varepsilon) = \frac{1}{2^{n(N-1)}} \sum_{\mathcal{B} \in \mathcal{C}(\varepsilon, N)} |B_\varepsilon^\cap(\mathcal{B})|^{-1}$$

with $B_\varepsilon^\cap(\mathcal{B}) = \bigcap_{i=1}^N B_\varepsilon(v_i)$. Lower and upper bounds for this probability are $V_\varepsilon^{N-1} \leq \mathbb{P}(\mathcal{B} \subset B_\varepsilon) \leq V_{2\varepsilon}^{N-1}$, where V_ε is the measure of an ε -ball.

Proof. First, in order to introduce the intuition to obtain the first result, observe that the proportion of template databases that can be covered by an ε -ball is bounded from above by counting the number of template databases included in each ball of \mathbb{Z}_2^n :

$$\mathbb{P}(\mathcal{B} \subset B_\varepsilon) \leq \frac{1}{2^{Nn}} \sum_{v \in \mathbb{Z}_2^n} |\mathcal{C}_v(\varepsilon, N)|$$

This quantity is an upper bound since each template database \mathcal{B} is counted several times, for different v . To be more specific, a template database is counted once for each $v \in B_\varepsilon^\cap(\mathcal{B})$. Next, notice that for a given database $\mathcal{B} \in \mathcal{C}_v(\varepsilon, N)$, $|B_\varepsilon^\cap(\mathcal{B})| \geq 1$, since v necessarily belongs to $B_\varepsilon^\cap(\mathcal{B})$. Then, the following equation can be written:

$$\mathbb{P}(\mathcal{B} \subset B_\varepsilon) = \frac{1}{2^{Nn}} \sum_{v \in \mathbb{Z}_2^n} \sum_{\mathcal{B} \in \mathcal{C}_v(\varepsilon, N)} |B_\varepsilon^\cap(\mathcal{B})|^{-1}.$$

Observe that the sum on $\mathcal{C}_v(\varepsilon, N)$ does not depend on v , so that:

$$\mathbb{P}(\mathcal{B} \subset B_\varepsilon) = \frac{1}{2^{Nn}} \sum_{v \in \mathbb{Z}_2^n} \sum_{B \in \mathcal{C}(\varepsilon, N)} |B_\varepsilon^\cap(B)|^{-1} = \frac{1}{2^{n(N-1)}} \sum_{B \in \mathcal{C}(\varepsilon, N)} |B_\varepsilon^\cap(B)|^{-1}$$

Next, to obtain the lower and upper bounds consider:

$$I_\varepsilon(v_1, \dots, v_k) = \cup_{v \in \cap_{i=1}^k B_\varepsilon(v_i)} B_\varepsilon(v)$$

and it follows that:

$$\mathbb{P}(\mathcal{B} \subset B_\varepsilon) = \mathbb{P}(v_2 \in B_{2\varepsilon}(v_1) \mid v_1) \times \dots \times \mathbb{P}(v_N \in I_\varepsilon(v_1, \dots, v_{N-1}) \mid (v_1, \dots, v_{N-1}) \subset B_\varepsilon)$$

Next, provide an upper bound for each previous term:

$$V_\varepsilon \leq \mathbb{P}(v_k \in I_\varepsilon(v_1, \dots, v_{k-1}) \mid (v_1, \dots, v_{k-1}) \subset B_\varepsilon) \leq V_{2\varepsilon}.$$

The upper bound is obtained by considering the case " $v_1 = \dots = v_{k-1}$ ". Concerning the lower bound, it is based on the fact that " $(v_1, \dots, v_{k-1}) \subset B_\varepsilon$ " implies that " $\cap_{i=1}^{k-1} B_\varepsilon(v_i)$ is a non-empty set, then $|I_\varepsilon(v_1, \dots, v_{k-1})| \geq |B_\varepsilon|$ ". ■

Theorem 2.2.4 gives the probability that there exists a template in \mathbb{Z}_2^n which impersonates all users if the threshold of ε is used. This result only refers to a rare event, but it is an intermediate result in order to provide Corollary 2.2.1, which focuses on the probability of a (k, ε) -master template. In the following, \mathcal{B}_k denotes a subset of k templates from the template database \mathcal{B} , and M_k is the event "an ε -ball covers k templates and none ball can not include these k templates plus another template from the template database".

Corollary 2.2.1. *The probability of a (k, ε) -master template is:*

$$\mathbb{P}(M_k) = \mathbb{P}(\mathcal{B}_k \subset B_\varepsilon) \times \frac{\binom{N}{k}}{2^{n(N-k)}} \left| \bigcap_{v \in B_\varepsilon^\cap(\mathcal{B}_k)} \overline{B_\varepsilon(v)} \right|^{N-k}$$

and lower and upper bounds for this probability are:

$$\binom{N}{k} V_\varepsilon^{k-1} (1 - V_{2\varepsilon})^{N-k} \leq \mathbb{P}(M_k) \leq \binom{N}{k} V_{2\varepsilon}^{k-1} (1 - V_\varepsilon)^{N-k}.$$

Proof. A (k, ε) -master template occurs when k templates can be covered with one (or more) ε -ball and the $N - k$ other templates do not belong to the covering ball(s):

$$\begin{aligned} \mathbb{P}(M_k) &= \binom{N}{k} \mathbb{P}(\mathcal{B}_k \subset B_\varepsilon) \mathbb{P}\left(\forall u \in \mathcal{B} \setminus \mathcal{B}_k, u \notin \bigcup_{v \in B_\varepsilon^\cap(\mathcal{B}_k)} B_\varepsilon(v) \mid \mathcal{B}_k \subset B_\varepsilon\right) \\ &= \binom{N}{k} \mathbb{P}((v_1, \dots, v_k) \subset B_\varepsilon) \prod_{j=1}^{N-k} \mathbb{P}\left(\bigcap_{v \in B_\varepsilon^\cap(\mathcal{B}_k)} \overline{B_\varepsilon(v)} \mid \mathcal{B}_k \subset B_\varepsilon\right) \\ &= \binom{N}{k} \mathbb{P}((v_1, \dots, v_k) \subset B_\varepsilon) \frac{1}{2^{n(N-k)}} \left| \bigcap_{v \in B_\varepsilon^\cap(\mathcal{B}_k)} \overline{B_\varepsilon(v)} \right|^{N-k} \end{aligned}$$

Next, in order to obtain lower and upper bounds, remark that, if $\mathcal{B}_k \subset B_\varepsilon$ holds, $B_\varepsilon^\cap(\mathcal{B}_k)$ is not empty, and that $\{v\} \subseteq B_\varepsilon^\cap(\mathcal{B}_k) \subseteq B_\varepsilon$, for an unknown specific $v \in \mathbb{Z}_2^n$. The lower case occurs if at least two templates u and $v \in \mathcal{B}_k$ are such that $d_H(u, v) = 2\varepsilon$ (they are opposed on an ε -sphere) and the upper case occurs when all templates of \mathcal{B}_k are equal. The results follow by replacing $B_\varepsilon^\cap(\mathcal{B}_k)$ with $\{v\}$ and then with B_ε . ■

To a lesser extent, the probability of a (k, ε) -near-collision with a given template $v \in \mathcal{B}$ is provided in Proposition 2.2.5 below. In other words, we consider a targeted version of a master-template, which is simpler since it is only required to evaluate the location with regards to this given template v , and not all the pairs of templates in the database. This result enables us to relate Corollary 2.2.1, which is about a complex event, to the case of an event that is simpler to understand from a computational point of view, since it leads to a Binomial distribution.

Proposition 2.2.5. *For a given template $v \in \mathcal{B}$, and $\mathcal{B}_{-v} = \mathcal{B} \setminus v$, the probability of a near-collision for v is:*

$$\mathbb{P}(\exists(v_1, \dots, v_k) \subset \mathcal{B}_{-v} \text{ such that } (v_1, \dots, v_k) \subset B_\varepsilon(v)) = \binom{N-1}{k} V_\varepsilon^k (1 - V_\varepsilon)^{N-k-1}$$

Proof. Recall that each $u \in \mathcal{B}_{-v}$ are independent and follows a uniform distribution on \mathbb{Z}_2^n . Denotes $v = (v_1, \dots, v_k)$ a vector in \mathcal{B}_{-v} , then $\mathbb{P}(\exists v \in \mathcal{B}_{-v} \text{ such that } v \subset B_\varepsilon(v))$ is equal to

$$\binom{N-1}{k} \prod_{v_j \in v} \mathbb{P}(v_j \in B_\varepsilon(v)) \prod_{v_\ell \notin v} \mathbb{P}(v_\ell \notin B_\varepsilon(v))$$

and the result follows with $\mathbb{P}(v_j \in B_\varepsilon(v)) = V_\varepsilon$. ■

2.2.5 Numerical Evaluations: Databases Security w.r.t Near-Collisions

The main problem associated with near-collisions is the drop in recognition performance for identification systems and the occurrence of master templates that facilitate multiple impersonations [30]. In the same experimental setup as the previous section, we can compute the critical population for reasonable scenarios with respect to Proposition 2.2.3. The critical population, denoted as N_r , is the maximum number of users such that the probability of a weak collision is smaller than a risk parameter r . In the literature on the birthday problem, it is classical to set $r = 1/2$. By modifying Proposition 2.2.3 and using the same method as the proof of the Proposition 2.2.3, it can be showed that given γ , n and ε the probability of a near-collision is smaller than a given r if

$$N_r \leq \frac{1}{2} + \sqrt{\frac{1}{4} - \frac{1}{48\gamma} - \frac{\log(1-r)}{V_\varepsilon}}.$$

Table 2.3 gives the critical population for scenarios of interest. As shown in this table, as long as the template size is greater than 128 and the threshold smaller than $n/10$, near-collisions should not be a problem for the actual human population *i.e.*, 10^{10} individuals.

2.2.6 Metric Based Analysis

The previous section is based on the analysis of the underlying metric space. The issue is that the methodology does not apply to all existing systems. To solve this issue, some authors in the

| n | 128 | | | 256 | | | 512 | | | 1024 | | | 2048 | | |
|------------------------|-----|----|----|-----|----|----|-----|----|-----|------|-----|-----|------|-----|-----|
| ε | 6 | 12 | 24 | 12 | 25 | 50 | 25 | 51 | 102 | 51 | 102 | 204 | 102 | 204 | 409 |
| $\log_{10}(N_{0.5})$ | 14 | 11 | 6 | 28 | 21 | 11 | 56 | 42 | 22 | 111 | 83 | 44 | 221 | 165 | 87 |
| $\log_{10}(N_{0.25})$ | 14 | 11 | 6 | 28 | 21 | 11 | 56 | 41 | 22 | 110 | 82 | 44 | 221 | 165 | 86 |
| $\log_{10}(N_{0.1})$ | 14 | 11 | 6 | 28 | 21 | 11 | 55 | 41 | 22 | 110 | 82 | 43 | 220 | 164 | 86 |
| $\log_{10}(N_{0.01})$ | 13 | 10 | 5 | 27 | 20 | 10 | 55 | 41 | 21 | 110 | 82 | 43 | 220 | 164 | 86 |
| $\log_{10}(N_{0.001})$ | 13 | 10 | 5 | 27 | 20 | 10 | 54 | 40 | 21 | 109 | 81 | 42 | 219 | 163 | 85 |

Table 2.3 – Critical population for various threshold and template sizes for $\gamma = 0.01$.

literature [127, 74, 70, 11] propose basing the analysis on measurements of biometric systems, such as FMR. The aforementioned analyses employ the birthday paradox to calculate the probability of a near-collision.

The Birthday Problem: The birthday problem in probability theory seeks to compute the likelihood of shared birthdays within a finite group of individuals. The probability $\mathbb{P}(n)$ of at least two people sharing a birthday within a group of size n randomly chosen can be calculated using the complement rule:

$$\mathbb{P}(n) = 1 - \frac{365 \cdot 364 \cdot \dots \cdot (365 - (n - 1))}{365^n}$$

More generally, in cryptography, we use the following setup. Let us consider a function with a value in a set E . If the output of the function is uniform in E , then the probability that in a set of n inputs, there are two equal outputs is:

$$\mathbb{P}(n) = 1 - \frac{|E| \times (|E| - 1) \times \dots \times (|E| - (n - 1))}{|E|^n}$$

The Biometric Birthday Problem: Daugman [127, 74, 70, 11] present the biometric birthday problem which is the biometric version of the birthday problem. Daugman states the biometric birthday problem as follows: “If some biometric technology is operating with a verification FMR, how many people, chosen at random, must be assembled until it becomes more likely than not that at least one pair of them have a biometric collision (are falsely matched to each other)?”

In its papers, Daugman claims that for N users, “if a biometric technology is operating at some verification False Match Rate FMR, then the probability of a given pairing not resulting in a False Match is $(1 - \text{FMR})$, and the probability that none of the possible pairings do so is $(1 - \text{FMR})^{N(N-1)/2}$ ”. The problem is that the interpretation of Daugman of the birthday problem is that for him, the probability $\mathbb{P}_D(N)$ of at least two people sharing a birthday within a group of size N randomly chosen is:

$$\mathbb{P}_D(N) = 1 - \frac{364^{N(N-1)/2}}{365}$$

This indicates that for $N = 366$, where we should have $\mathbb{P}(N) = 0$, we have $\mathbb{P}_D(N) > 0$. The error in this reasoning is the assumption of independence, which is not valid in this context. Indeed, if the pairing of A and B matches and the pairing of B and C matches, then, with a high probability, the pairing of A and C will match. Consequently, the information regarding previous pairings affects the probability of other pairings matching. The main problem with this result is that considering

the independence gives an approximation and not the real value [21].

Details on Daugman Biometric Birthday Problem Given the lack of details in these papers, we provide the full analysis to derive this approximation. First, a precision on the FMR. The FMR is fixed for a given biometric system in the authentication mode and does not change even if we increase the database population. The value that could change is the estimation of the FMR but in the following, we suppose that we have access to the theoretical False Match Rate denoted by FMR. For more details on the FMR, please refer to Section 1.3.2. Let M_i be the event "The individuals within the i -th pair of $\bar{\mathcal{P}}$ does not falsely match", $Q(N)$ the probability of the event "There is no false match among N pairs", and $P(N)$ the probability of the event "There is at least one false match within N pairs". We have $\mathbb{P}(M_i) = \text{FMR}$ and $\mathbb{P}(K) = 1 - Q(K)$. By denoting $\bar{\mathcal{P}}$ be the set of all possible unordered pairs among N individuals then, the number of all possible pairs of individuals is given by

$$|\bar{\mathcal{P}}| = \binom{N}{2} = \frac{N!}{2!(N-2)!} = \frac{N(N-1)}{2}.$$

Then, we have:

$$\begin{aligned} Q(N) &= \mathbb{P}\left(\bigcap_{i=1}^{|\bar{\mathcal{P}}|} M_i\right) \\ &= \mathbb{P}(M_1) \times \prod_{i=2}^{|\bar{\mathcal{P}}|} \mathbb{P}\left(M_i \mid \bigcap_{j=1}^{i-1} M_j\right) \\ &\approx \mathbb{P}(M_1) \times \prod_{i=2}^{|\bar{\mathcal{P}}|} \mathbb{P}(M_i) \end{aligned}$$

The above approximation is a key part of the result from which we can deduce:

$$Q(N) \approx (1 - \text{FMR})^{\frac{N(N-1)}{2}}$$

by using the fact that $\forall i \in \{1, \dots, |\bar{\mathcal{P}}|\}$, $\mathbb{P}(M_i) = 1 - \text{FMR}$. Then, we have the following result:

$$P(K) \approx 1 - (1 - \text{FMR})^{\frac{N(N-1)}{2}}$$

which explains the result of Daugman.

Critical Population Size: To achieve a given level of security and accuracy for the occurrence of near-collision, the number of clients in a database must be bounded. We say that the population of a database is critical when the probability of a near-collision is greater than $1/\lambda$ with $\lambda \in \mathbb{R}_{\geq 2}$. Daugman [11] gives an estimation of the critical population for $\lambda = 2$ which is $N \approx \sqrt{\frac{1.386}{\text{FMR}}}$. We generalize it to any $\lambda \in \mathbb{R}_{\geq 2}$.

Theorem 2.2.5. *Given a biometric system operating at some verification False Match Rate FMR then, an*

| FMR | System | Modality | Critical Population |
|-------------------------|---------------------|------------------|---------------------|
| 1 in 5.0×10^5 | Apple Touch ID [4] | Fingerprint | 833 |
| 1 in 1.0×10^6 | Apple Face ID [3] | Face | 1,178 |
| 1 in 2.0×10^5 | Google standard [7] | All | 527 |
| 1 in 1.0×10^5 | NEC5 [58] | Iris | 373 |
| 1 in 1.0×10^4 | Nikisins [57] | VBR | 118 |
| 4 in 1.0×10^2 | ASEC [66] | Online signature | 6 |
| 11 in 1.0×10^4 | NXOR [66] | Face | 36 |
| 5 in 1.0×10^5 | HXKJ [6] | Fingerprint | 167 |
| 3 in 1.0×10^4 | MM_PV [6] | Palm vein | 68 |
| 15 in 1.0×10^4 | Biotope [6] | STF | 31 |

Table 2.4 – Critical population comparison in function of FMR with $\lambda = 2$.

estimation of the critical population for a given λ is given by:

$$N \approx \left\lfloor \frac{1}{2} + \sqrt{\frac{1 + \frac{8 \ln(1-\lambda^{-1})}{\ln(1-\text{FMR})}}{4}} \right\rfloor$$

Proof. Given that the probability that there is a near-collision is approximately $P(N) = 1 - (1 - \text{FMR})^{\frac{N(N-1)}{2}}$, we can infer the approximation critical population size. We seek N such that

$$\begin{aligned} 1 - (1 - \text{FMR})^{\frac{N(N-1)}{2}} &\geq 1/\lambda \\ \frac{N(N-1)}{2} \ln(1 - \text{FMR}) &\leq \ln(1 - \lambda^{-1}) \\ N^2 - N - 2 \frac{\ln(1 - \lambda^{-1})}{\ln(1 - \text{FMR})} &\leq 0. \end{aligned}$$

The study of the function yields the approximate critical size of the population:

$$N \leq \frac{1}{2} + \sqrt{\frac{1 + \frac{8 \ln(1-\lambda^{-1})}{\ln(1-\text{FMR})}}{4}}.$$

■

Remark 2.2.2. If we take $\lambda = 2$ in Theorem 2.2.5, the critical population is approximately what Daugman has for $\lambda = 2$.

Table 2.4 gives critical population for $\lambda = 2$ for several systems. We can see that the critical population is very low for the majority of the systems.

We also provide an estimation for the FMR required to manage a given of size population N , with respect to λ .

Corollary 2.2.2. Let N be the population of a biometric database and $\lambda \in \mathbb{R}_{\geq 2}$. Then, an estimation for the maximal FMR such that the probability of a near collision is smaller than $1/\lambda$ is given by:

$$E_{\text{FMR}} = 1 - e^{-2 \frac{\ln(1-1/\lambda)}{N(N-1)}}$$

| | | | | | | | | | |
|------------------------------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Population Size (\log_{10}) | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| FMR Needed (\log_2) | -13 | -19 | -26 | -33 | -39 | -46 | -53 | -59 | -66 |
| Security bits Needed | 13 | 20 | 27 | 33 | 40 | 46 | 53 | 60 | 66 |

Table 2.5 – FMR maximal needed in function of the population size of the system for $\lambda = 2$.

Proof. Similarly to the previous proof, we want:

$$1 - (1 - \text{FMR})^{\frac{N(N-1)}{2}} \geq 1/\lambda$$

$$\ln(1 - \text{FMR}) \leq 2 \frac{\ln(1 - \lambda^{-1})}{N(N-1)}$$

and the result follows. ■

The FMR needed to ensure that the probability of a near collision is smaller than $1/2$ for the world population ($N = 10^9$) is $\text{FMR} \approx 2^{-66}$. Table 2.5 gives the FMR for interesting population size.

2.3 Targeted and Untargeted Attack Models

The objective of the attacker in the two subsequent sections is to either recover the stored biometric data $x \in \mathbb{Z}_q^{n-1}$ with $q \geq 2$ or to impersonate a user. This represents one of the most devastating types of attack, as described in Section 1.3.4. The attacker does not have access to the database. To emulate the interaction between the biometric system and the attacker, he has access to an oracle, designated as $\mathcal{O}_{x,\varepsilon}$. The oracle received the template selected by the attacker and compared it with the template that had been previously enrolled and stored. If the distance is below the threshold ε , the oracle returns 1 and 0 otherwise. In a more formal way, $\mathcal{O}_{x,\varepsilon}$ is a function defined as:

$$\mathcal{O}_{x,\varepsilon} : \mathbb{Z}_q^n \rightarrow \{0, 1\}$$

$$y \mapsto \begin{cases} 1 & \text{if } d_H(x, y) \leq \varepsilon. \\ 0 & \text{otherwise.} \end{cases}$$

Templates are analogous to hashed passwords. The objective is to conduct a similar analysis of the biometric template as that of the well-known analysis of hashed passwords. Firstly, targeted attacks are considered. The objective is to identify the stored biometric data of a specific user or a template that can impersonate this user. This is analogous to finding a collision for a specific hashed password. Secondly, untargeted attacks are considered. In this setting, the adversary's objective is to impersonate a user within a database without focusing on any particular individual. In the context of these attacks, it is assumed that the attacker has access to the attack Point 4 for the minimal setup and gains control of Point 8 in the leakage scenarios. For a more comprehensive understanding of these threat points, please refer to Figure 1.2 and Section 1.3.3. Figure 2.7 shows the main difference between targeted and untargeted attacks.

1. Most of our results are illustrated in the case $q = 2$.

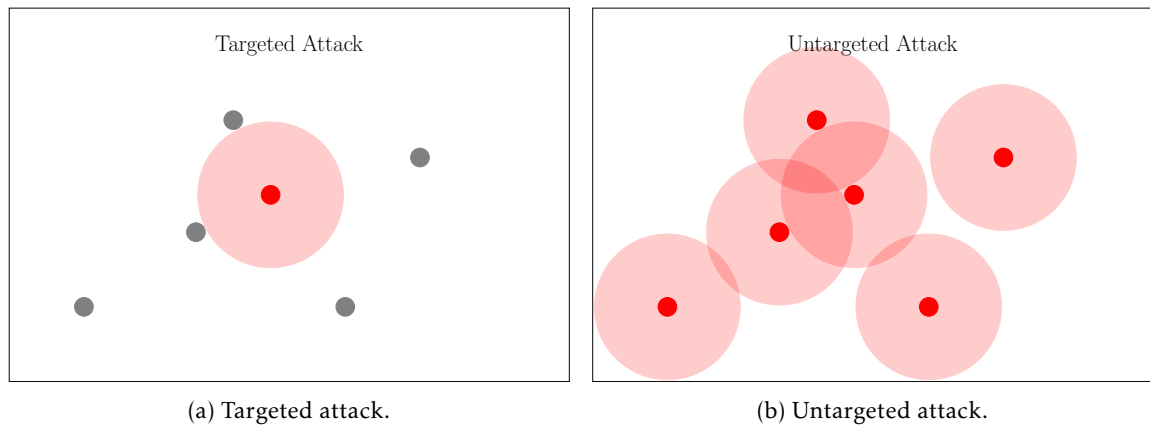


Figure 2.7 – Difference between targeted and untargeted attacks from the attacker point of view. The highlighted surface represents the attackable area.

To model the behavior of an ideal biometric system, we assume that the templates are uniformly distributed across their respective space. This assumption is pivotal as it enables us to establish absolute upper bounds that apply regardless of the actual distribution of the templates. In practice, if the distribution is non-uniform, some regions of the space will have a higher concentration of templates. An attacker who is aware of these variations in template density could exploit this knowledge to design a more effective attack strategy, rather than assuming a uniform distribution. Thus, while the uniform distribution assumption may be idealized, it provides a useful benchmark for evaluating the performance and security of biometric systems.

2.4 Targeted Attack

This section is devoted to the analysis of targeted attacks using exhaustive search methods in the context of minimal leakage scenarios. Subsequently, the results are expanded to encompass different leakage scenarios. In this setup, the attacker does not have direct access to the database and interacts with an oracle, as defined in Section 2.3. The objective of the attacker is to retrieve the enrolled template in the least number of calls to the oracle.

2.4.1 Exhaustive Search

In this section, the attacker aims to find a template that lies in the ball of center x (the target template) and radius ε (the threshold). To identify such a point, several methods are available, each with its own set of advantages and disadvantages. In this section, in order to provide the most general result, we assume that the metric space is \mathbb{Z}_q^n equipped with the Hamming distance d_H .

Brute Force: The objective of this attack is to exhaustively test all potential points in space until the oracle $\mathcal{O}_{x,\varepsilon}$ yields 1. In the worst case, we test every point in space, which results in the examination of q^n vectors. To obtain this result, we ignore the ε acceptance threshold. On the other hand, if we consider that only $n - \varepsilon$ exact coordinates are needed to be accepted by the system, complexity decreases to $q^{n-\varepsilon}$ tests. Since the attacker specifically targets $n - \varepsilon$ coordinates (the attacker arbitrarily chooses ε coordinates that do not change), and aims for a perfect match for the $n - \varepsilon$ remaining coordinates yielding the result.

Random Sampling: The attacker randomly chooses a template in \mathbb{Z}_q^n and tests it by querying the oracle $\mathcal{O}_{x,\varepsilon}$. The worst case for the attacker is that the templates are uniformly distributed in \mathbb{Z}_q^n . The probability that a template submitted to $\mathcal{O}_{x,\varepsilon}$ yields 1 is $\rho = \frac{|B_{q,\varepsilon}(x)|}{q^n}$. According to this naive strategy, we can assume that the tests are independent and that each can be represented by a Bernoulli experiment with a success probability of ρ . The number of tries needed before one success follows a geometric distribution. Then, the expected number of tries for an attacker to get accepted by the system is ρ^{-1} . First, recall that the cardinal of $B_{q,\varepsilon}(x)$ is

$$|B_{q,\varepsilon}(x)| = \sum_{i=0}^{\varepsilon} \binom{n}{i} (q-1)^i,$$

and that the q -ary entropy is $h_q(x) = x \log_q(q-1) - x \log_q x - (1-x) \log_q(1-x)$. Then, using the Stirling approximation (see [72, 119]), the expected number of tries for an attacker is

$$\rho^{-1} = \frac{q^n}{|B_{q,\varepsilon}(x)|} = \frac{q^n}{\sum_{i=0}^{\varepsilon} \binom{n}{i} (q-1)^i} \leq \frac{q^n}{q^{nh_q(\varepsilon/n)+o(n)}} = q^{n(1-h_q(\varepsilon/n))+o(n)}$$

if $\frac{\varepsilon}{n} \leq 1 - \frac{1}{q}$ holds and if n is large enough.

Random Sampling Without Point Replacement: As the random sampling, the attacker randomly chooses a template in the set $S \subseteq \mathbb{Z}_q^n$. At each step, if $\mathcal{O}_{x,\varepsilon}$ returns 0, the tested vector b is removed from the set S . The probability of success does not remain constant throughout the experiment, unlike in the previous case. Consequently, the experiment follows a hypergeometric distribution. This game is equivalent to having an urn with q^n object where $|B_{q,\varepsilon}(x)|$ are considered "good". Then, according to Ahlgren [79] the expected number of queries to $\mathcal{O}_{x,\varepsilon}$ before success is given by

$$\frac{q^n + 1}{|B_{q,\varepsilon}(x)| + 1} \approx \rho^{-1}.$$

This attack has a slightly better performance compared to the previous one, although it is accompanied by an exponential memory cost that reduces its efficiency, making this version less interesting than the previous one.

Remark 2.4.1. *In the case of random sampling, if the value of n is large, it is preferable to select a draw with replacement to save memory while maintaining a high degree of performance. Indeed, the probability of drawing a vector that has already been selected is relatively small if n is sufficiently large.*

Tree Search: This algorithm was proposed by Pagnin *et al.* [83]. The underlying idea is to construct a tree of depth n such that each point of the space is considered to be a leaf. The tree structure is utilized to establish relative relations among the points of \mathbb{Z}_q^n and to guarantee that after each unsuccessful trial, non-overlapping portions of the space \mathbb{Z}_q^n can be removed. Specifically, if a point $p \in \mathbb{Z}_q^n$ does not satisfy the authentication, the algorithm removes not only the tested point p from the set of potential centers but also its sibling relatives generated by the common ancestor ε . At each try, the attacker can remove approximately q^ε points from the research space (for more details, please refer to [83]). The running time of the attack is the cost of exploring a q -ary tree of order

$n - \varepsilon$.

Remark 2.4.2. *It should be noted that as intended, the cost of all the presented attacks is exponential.*

Optimal Solution: The optimal solution is to solve the set-covering problem [94] using balls of radius ε . The main idea is to cover the space with the smallest number of balls of radius ε to partition the space. The objective is to remove an entire ball of radius ε if the query fails. This is an instance of the set covering problems. Pagnin *et al.* [83] claimed that the number of points that the adversary needs to query is only a factor of $O(\varepsilon \ln(n+1))$ more than the optimal cover. However, the result is imprecise, as detailed below in this remark, mainly because the optimal cover is not given. The strategy between a bounded and an unbounded adversary may differ as detailed in the following.

- **Unbounded adversary:** The adversary solves the NP-hard set covering problem [94] to find the optimal covering of \mathbb{Z}_q^n using balls of radius ε . The adversary exhaustively searches x using at most $q^{n(1-h_q(\varepsilon/n))+o(n)}$ queries to $\mathcal{O}_{x,\varepsilon}$. The number of vectors involved in a given optimal cover is $\frac{q^n}{|B_{q,\varepsilon}(x)|}$, which can be asymptotically approximated as detailed in what follows. Then, using bounds on the binomial coefficient (see [119, 72]), the result follows if $\frac{\varepsilon}{n} \leq 1 - \frac{1}{q}$ holds and if n is large enough.
- **Bounded adversary:** The adversary may use a greedy algorithm to find a non-optimal covering containing $\frac{q^n H(n)}{|B_{q,\varepsilon}|}$ vectors [143] with $H(n) = \sum_{i=1}^n i^{-1}$ the n -th harmonic number. The adversary then finds a solution with an exhaustive search in at most $\frac{q^n H(n)}{|B_{q,\varepsilon}|}$ queries. In order to provide a more intuitive value, notice that $\frac{q^n H(n)}{|B_{q,\varepsilon}|}$ can be bounded up by $\frac{q^n (\ln(n)+1)}{|B_{q,\varepsilon}|}$. As in the unbounded case, using the q -ary entropy and Stirling's approximation, this non-optimal covering leads the attacker to make at most $q^{n(1-h_q(\varepsilon/n))+o(n)}$ queries, as $\log_q(\ln(n)+1) = o(n)$.

Remark 2.4.3. *The time required to configure the greedy algorithm is exponential, rendering the aforementioned attack impractical. Moreover, even if an attacker computes the optimal covering, it still needs to query an exponential number of time the $\mathcal{O}_{x,\varepsilon}$ to find a point close to x .*

It is also interesting to note that the expected time for an attacker to be accepted by the system ($\mathcal{O}_{x,\varepsilon}$ gives 1) using the random sampling with and without replacement method is equivalent to the worst case using the optimal method.

Example of Expectations for the Random Sampling: To illustrate the influence of the threshold and the choice of q on exhaustive search, we calculate the precise expectation of the number of attempts required for an attacker to authenticate in different settings using the random sampling method. The results are presented in Table 2.6a and Table 2.6b. Experimental results show that to increase the security against exhaustive search, it is more interesting to increase q than to decrease ε .

2.4.2 The Center Search Attack

This section is based on the work of Pagnin *et al.* [83], to which we have added experiments, clarifications, and expanded the results. The possibility of a center search attack arises when the server requires only a template from the client in order to authenticate it. This may be the case when either point 3 or 4 has been compromised, as illustrated in Figure 1.2 in Section 1.3.3. The attack described in [83], consists of 3 distinct phases. These are exhaustive search (Section 2.4.1),

| System | n | ε | q | RSR (\log_2) | | n | ε | q | RSR (\log_2) |
|-------------------|-------|---------------|-----|---------------------|--|-------|---------------|-----|---------------------|
| IrisCode [104] | 2,048 | 738 | 2 | 121.37 | | 2,048 | 1,024 | 2 | 0.97 |
| IrisCode [104] | 2,048 | 656 | 2 | 199.94 | | 2,048 | 512 | 2 | 391.54 |
| IrisCode [104] | 2,048 | 574 | 2 | 300.24 | | 2,048 | 256 | 2 | 939.79 |
| FingerCode [15] | 80 | 30 | 2 | 5.92 | | 2,048 | 1,024 | 3 | 178.83 |
| BioHashing [84] | 180 | 60 | 2 | 17.74 | | 2,048 | 512 | 3 | 1,077.86 |
| BioEncoding [100] | 350 | 87 | 2 | 70.62 | | 2,048 | 256 | 3 | 1,881.91 |
| BioEncoding [100] | 350 | 105 | 2 | 45.18 | | 2,048 | 1,024 | 4 | 430.24 |

(a) Examples with real biometric systems.

(b) Examples with several q and ε .

Table 2.6 – Expected number of calls to oracle for the exhaustive search method Random Sampling with Replacement (RSR).

edge detection, and center search. In this section, we assume that we have a point y which is within the ball centered in x and with radius ε . Pagnin *et al.* [83] shows that to determine the center of the ball $B_{2,\varepsilon}(x)$ in the binary case, we proceed in two stages:

1. Find the Edge of the ball.
2. Find the center using a Hill Climbing attack.

Edge Detection: To identify the edge of the ball, the coordinates of y are modified one by one and each new vector y' is tested. If y' is rejected, then, y is on the edge of the ball. If y' is within the ball, we state $y = y'$ and we modify the next coordinate until the template is no longer within the ball. The complexity of this step is at most 2ε which corresponds to the case where y was already on an edge and we cross the ball $B_{2,\varepsilon}(x)$ by following its diameter. This algorithm is linear in ε .

Center Search: Let us consider the case where y is a point lying on the edge of the ball $B_{2,\varepsilon}(x)$ for the binary case. The idea of the method is the following. As y lies at the edge of the ball, if we change a coordinate and stay within the ball then, the new point is closer to the center. If we apply this method to all n coordinates, we completely determine the center of the ball x . The cost of this part is linear in n (exactly n queries).

Generalize the Center Search Attack: Using the same reasoning, the center search attack is trivially generalized to \mathbb{Z}_q^n equipped with the Hamming distance. The cost of the edge detection does not change but the cost of the center search goes from n to $n(q-1)$. Algorithm 2 provides a constructive proof of the generalization of the center search for any $q \geq 2$.

2.4.3 The Impact of Information Leakage on Targeted Attack

In the previous section, we examined the case of minimal leakage. That is, when comparing rich and enrolled data, the entity in charge of the comparison only returned 0 (does not match) or 1 (match). However, in the literature, we find primitives that leak more information under certain conditions. In this section, we study the impact of these information leaks on the attacker's naive strategy, in order to provide complexities that can be applied to all systems without worrying about these specificities. As in the previous section, the attacker makes queries to $\mathcal{O}_{x,\varepsilon}$ to know if its template matches or not the enrolled one. In this case, $\mathcal{O}_{x,\varepsilon}$ provides additional information as detailed later in this manuscript.

Algorithm 2 : Center Search for any $q \geq 2$

Input : $y \in \mathbb{Z}_q^n$ a vector lying on the edge of the ball $B_{q,\varepsilon}(x)$.
Output : $x \in \mathbb{Z}_q^n$ the center of the ball $B_{q,\varepsilon}(x)$.

```

1 Function Center_Search( $y$ ):
2    $x$  is initialized to be the same as  $y$ ;
3   for  $j \leftarrow 1$  to  $n$  do
4      $v$  is initialized to be the same as  $y$ ;
5     for  $l \leftarrow 1$  to  $q - 1$  do
6       The  $j$ -th coordinate of  $v$  is incremented by 1 modulo  $q$ ;
7       if The query of  $v$  to  $\mathcal{O}_{x,\varepsilon}$  yields 1 then
8         Set the  $j$ -th coordinate of  $x$  to the  $j$ -th coordinate of  $v$  ;
9       end
10    end
11  end
12 return  $x$ ;
```

| Distance-to-Threshold comparison | Leakage | Complexity type | Complexity in Big-O | Theorem |
|----------------------------------|-------------------------------------|-------------------------|---|---------|
| Below | Distance | Exponential | $q^{n-\varepsilon} + q\varepsilon$ | 2.4.1 |
| | Positions | Exponential | $q^{n-\varepsilon} + q$ | 2.4.2 |
| | Positions and values | Exponential | $q^{n-\varepsilon}$ | 2.4.3 |
| | Positions and values (accumulation) | Linearithmic/Polynomial | $n^\alpha \log n$ | 2.4.7 |
| Both | Minimal ¹ | Exponential | $q^{n-\varepsilon} + n(q-1) + 2\varepsilon$ | [83] |
| | Distance | Linear | nq | 2.4.4 |
| | Positions | Constant | q | 2.4.5 |
| | Positions and values | Constant | 1 | 2.4.6 |

Table 2.7 – Summary of all leakage exploits and their complexities with α such that the occurrence of the rarest error is $n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$. The Distance-to-Threshold comparison determines if the leak occurs when $d(x, y) \leq \varepsilon$ (below) or when there is no distance requirement between x and y (both). For all the complexities, x and y are in \mathbb{Z}_q^n with $q \geq 2$ except for the minimal leakage where x and y are in \mathbb{Z}_2^n . The provided complexities represent worst-case scenarios, except for the accumulation attack where the result is the expectation.

¹Note that the Big-O complexity of the optimal exhaustive search strategy, in the worst-case, is the same as the naive strategy as the minimum of $h(\cdot)$ is 0.

The contributions consist of various information leakage scenarios, the corresponding generic attacks, and their complexities. The discussed scenarios give rise to new attacks:

- Accumulation attacks that capture potential attacks from an *honest-but-curious* server during a client authentication. These attacks assume the use of a privacy-preserving cryptographic scheme for evaluating the distance between two hidden templates. Specifically, as an example, we assume the use of a cryptographic obfuscator for the distance function.
- Attacks from malicious clients exploiting various information leakages from the matcher in the context of a leaked (but obfuscated) database, or by interacting with the server during an *online exhaustive search attack*.

The complexities of the attacks, relying on different scenarios, are summarized in Table 2.7.

Example of Information Leakage

A matcher built on top of an inner product functional encryption with a function privacy scheme whose functionality is to provide a secret key for evaluating an inner product [43, 44]. The desired distance function (*e.g.*, Hamming distance) can be formulated as an inner product as shown in Section 2.2.2. The security is given in adversarial models that rule out simple attacks like

hill-climbing attacks, as the scheme directly leaks the distance.

Typology of Information Leakage

Except for the accumulation attack, the attacker exploits points 4 and 8 in all discussed scenarios. Point 4 allows the submission of a chosen template, while Point 8 grants access to additional information beyond the binary output. The accumulation attack only necessitates control over the Point 8. For detailed insights into the remaining threat points, readers are referred to Figure 1.2 and Section 1.3.3. There are three main categories of information leakage:

- Below the threshold.
- Above the threshold.
- Both below and above the threshold.

In each of these categories, several sub settings can be identified. The first one corresponds to the absence of any leakage, resulting in $\mathcal{O}_{x,\varepsilon}$ yielding only the binary output. Then, the following information leakages are examined:

- The distance.
- The positions of the errors.
- Both the error positions and values.
- Both the distance and the positions of the errors.
- Both the distance and the positions and their corresponding erroneous values.

It is not relevant to consider that additional information is leaked only above the threshold, as no scheme has such behavior. As a consequence, solely scenarios ‘below the threshold’ and ‘below and above the threshold’ are examined. The Hamming distance is a measure of the number of differing coordinates between two templates. Therefore, knowledge of the erroneous coordinates implies knowledge of the distance itself. Hence, we do not consider all possible scenarios.

2.4.4 Exploiting the Leakage

Attack Complexities for Leakage Below the Threshold

Leakage below the threshold is considered in this section. Given the hidden target x , querying y such that $d_H(x, y) \leq \varepsilon$ to the oracle $\mathcal{O}_{x,\varepsilon}$ provides information beyond the binary output. Concerning fuzzy matchers that employ secure sketches or error-correcting code mechanisms, information leakage deliberately occurs below the threshold for error-correction purposes. In many instances, information related to the distance, error locations, and errors themselves are either explicitly calculated or can be inferred.

Leakage of the Distance: The first case occurs when the distance is given to the attacker as extra information.

Theorem 2.4.1. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the distance below the threshold, an attacker can retrieve x in the worst case in $O(q^{n-\varepsilon} + q\varepsilon)$ queries to $\mathcal{O}_{x,\varepsilon}$.*

Proof. The matcher, using the Hamming distance, requires a minimum of $n - \varepsilon$ accurate coordinates to output 0. Since the attacker specifically targets $n - \varepsilon$ coordinates (the attacker arbitrarily chooses ε coordinates that do not change), an exhaustive search attack is performed in at most $q^{n-\varepsilon}$ steps to get accepted by the matcher. Then, a hill-climbing attack runs on the remaining ε coordinates to minimize the distance at each step. Coordinate by coordinate, the attacker obtains the right value if

the distance decreases. Since there are q different values to test on ε coordinates, determining the correct ones requires a maximum of $(q-1)\varepsilon$ steps. Then, the overall complexity is $O(q^{n-\varepsilon} + q\varepsilon)$. ■

Leakage of the Positions: The positions of the errors are the extra information given to the attacker, while their values remain secret.

Theorem 2.4.2. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the positions of the errors below the threshold, an attacker can retrieve x in the worst case in $O(q^{n-\varepsilon} + q)$ queries to $\mathcal{O}_{x,\varepsilon}$.*

Proof. As the leakage occurs solely below the threshold, the first step is to find a vector $y \in \mathbb{Z}_q^n$ such that $d(x, y) \leq \varepsilon$. To identify such a vector, the attacker performs an exhaustive search attack in $q^{n-\varepsilon}$ steps, as previously shown. Since ε coordinates remain unknown, and each coordinate ranges from 0 to $q-1$, every possibility must be examined. By testing all possibilities simultaneously – for instance, testing all coordinates at 0, then all coordinates at 1, and so forth up to $q-2$ while retaining the correct values – the original vector can be identified in no more than $q-1$ queries (refer to the example illustrated in Figure 2.8). Therefore, the complexity of the attack for recovering x is $O(q^{n-\varepsilon} + q)$. ■

Figure 2.8 gives a representation of the attack described above in the case \mathbb{Z}_4^5 and the hidden vector or the missing coordinates is $(0, 1, 3, 2, 2)$. Note that the actual complexity is $q-1$ since the final exchange is unnecessary, as the coordinates at $q-1$ become known after $q-1$ queries by inference.

Leakage of the Positions and the Values: When a vector below the threshold is given to the oracle $\mathcal{O}_{x,\varepsilon}$, the attacker gets information about both error positions and their values. This is similar to an error-correction mechanism designed to correct errors below a given threshold. Note that in the binary case, this scenario is the same as the previous one, hence the only considered case is $q > 2$.

Theorem 2.4.3. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the positions and the values of the errors below the threshold, an attacker can retrieve x in $O(q^{n-\varepsilon})$ queries to $\mathcal{O}_{x,\varepsilon}$.*

Proof. First, an exhaustive search is performed to find a vector y for which the distance is below the threshold, for a cost of $O(q^{n-\varepsilon})$. Then, given the error positions and the corresponding error values, y yields immediately the recovery of x . In the end, the complexity of the attack is $O(q^{n-\varepsilon})$. ■

Leakage Below and Above the Threshold

The second scenario is considered in this section, which involves a leakage independent of the threshold. In other words, when a hidden vector x is targeted, the queried vector y to the oracle $\mathcal{O}_{x,\varepsilon}$ results in the leak of additional information.

Leakage of the Distance: In this case, $d(x, y)$ the distance between $y \in \mathbb{Z}_q^n$ the fresh template and $x \in \mathbb{Z}_q^n$ the old template is leaked to the attacker regardless of the threshold.

Theorem 2.4.4. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the distance, an attacker can retrieve x in $O(nq)$ queries to $\mathcal{O}_{x,\varepsilon}$.*

| | | | | | |
|------------------|---------------|--------------|--------------|--------------|--------------|
| <u>Queries:</u> | $(\boxed{0},$ | $0,$ | $0,$ | $0,$ | $0)$ |
| | $(1,$ | $\boxed{1},$ | $1,$ | $1,$ | $1)$ |
| | $(2,$ | $2,$ | $2,$ | $\boxed{2},$ | $\boxed{2})$ |
| | $(3,$ | $3,$ | $\boxed{3},$ | $3,$ | $3)$ |
| <u>Solution:</u> | $(0,$ | $1,$ | $3,$ | $2,$ | $2)$ |

Figure 2.8 – Exploiting the error position leaked in the case \mathbb{Z}_4^5 and the hidden vector or missing coordinates is $(0, 1, 3, 2, 2)$.

Proof. As the attacker has access to the distance, it is possible to perform a hill-climbing attack, trying to minimize the distance at each step. The strategy is to find the vector y , coordinate by coordinate. As each coordinate has q possible values and there are n coordinates, this is done in $O(nq)$ steps. ■

Leakage of the Positions: The extra information given to the attacker is the positions of the errors.

Theorem 2.4.5. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the positions of the errors, an attacker can retrieve x in $O(q)$ queries to $\mathcal{O}_{x,\varepsilon}$.*

Proof. She tries the vector $(0, \dots, 0)$, $(1, \dots, 1)$ up to, $(q-1, \dots, q-1)$ and keep for each coordinate the right value (see Figure 2.8). Hence, the complexity of the attack to recover x is $O(q)$. ■

Leakage of the Positions and the Values: In this last case, the positions of the errors and corresponding values are leaked. Unlike the scenario of leakage below the threshold, such a leak provides an error-correcting code mechanism that operates irrespective of any distance and threshold.

Theorem 2.4.6. *Given ε a threshold, $x \in \mathbb{Z}_q^n$ a vector, and $\mathcal{O}_{x,\varepsilon}$ providing the positions of the errors and their values, an attacker can retrieve x in $O(1)$ queries to $\mathcal{O}_{x,\varepsilon}$.*

Proof. The submission of any vector gives the position of each error, and how to correct them, yielding a complexity in $O(1)$. ■

2.4.5 Accumulation Attack: A Passive Attack

During the client authentications, the attacker passively gathers information by observing errors leaked by the server. More specifically, the server leaks a list of positions and errors computed over the integers (*i.e.*, $x_i - y_i$) made by a genuine client during each authentication. Such information gathered during one successful authentication attempt is called an observation. The attacker aims to partially or fully reconstruct x by exploiting these observations.

In the binary case (*i.e.*, $q = 2$), the errors precisely yield the bits. If $x_i - y_i = 1$ then $x_i = 1$, and if $x_i - y_i = -1$ then $x_i = 0$. This attack is related to the Coupon Collector's problem [81], which involves determining the expected number of rounds required to collect a complete set of distinct coupons, with one coupon obtained at each round, and each coupon acquired with equal probability.

Exemple 2.4.5.1. *Suppose a setting with a metric space \mathbb{Z}_2^n equipped with the Hamming distance. A client seeks to authenticate to an honest-but-curious server that uses a scheme leaking $d(x, y)$ and the*

corresponding errors if $d(x, y) \leq \varepsilon$. As the client is legitimate, i.e., $d(x, y) \leq \varepsilon$ with a high probability, the attacker recovers the values of at most ε erroneous bits. The attacker needs to collect all the bits of the client, turning this problem into a Coupon Collector problem. For example, let assume $x = (0, 0, 1, 1, 0, 1, 0)$, $\varepsilon = 3$. The attacker sets $z = (?, ?, ?, ?, ?, ?, ?)$. Session 1: The client authenticates with $y = (1, 1, 0, 1, 0, 1, 0)$. In this case, $d(x, y) = 3 \leq \varepsilon$. The values of the erroneous bits of the client are obtained, yielding $z = (0, 0, 1, ?, ?, ?, ?)$. Session 2: the client authenticates with $y = (0, 0, 0, 0, 1, 1, 0)$. In this case, $d(x, y) = 3 \leq \varepsilon$, and the attacker obtains the value of the erroneous bits of the client and updates $z = (0, 0, 1, 1, 0, ?, ?)$. At this point, replacing the unknown values with random bits gives a vector that lies inside the acceptance ball as the number of unknown coordinates is smaller than the threshold ε .

In biometrics, some errors happen more frequently than others. In this setup, the Weighted Coupon Collector's Problem must be considered. Each coupon (i.e., each error) has a probability p_i to occur. Suppose that $p_1 \leq p_2 \leq \dots \leq p_n$ and $\sum_{i=1}^n p_i \leq 1$ then, according to Berenbrink and Sauerwald [103] (Lemma 3.2), the expected number of round E is such that:

$$\frac{1}{p_1} \leq E \leq \frac{H(n)}{p_1}$$

with $H(n)$ the n -th harmonic number. The upper bound on $H(n)$ is $1 + \log n$, which yields the expected number of rounds required to complete the collection:

$$\frac{1}{p_1} \leq E \leq \frac{\ln(n) + 1}{p_1}.$$

However, while in the original problem one coupon is obtained at each round, the number of errors made by a client during an authentication session is variable, i.e., between 1 and ε . In this case, the expected number of rounds required before all the errors have been observed is smaller than in the case where only one error occurs at each round. Consequently, the expected number of rounds required to collect all the errors is still in $\mathcal{O}(\log n/p_1)$.

Theorem 2.4.7. *Given ε a threshold, $x \in \mathbb{Z}_2^n$ a vector, $\text{Match}_{x,\varepsilon}$ leaks the positions of the errors and their values below the threshold, and assuming that the rarest coupon is obtained with probability $p_1 = n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$ an attacker can retrieve x in $\mathcal{O}(n^\alpha \log n)$.*

Proof. According to the Weighted Coupon Collector's problem and assuming that the rarest coupon is obtained with probability $p_1 = n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$, the vector x is recovered in $\mathcal{O}(n^\alpha \log n)$ observations. ■

It is worth noting that in this scenario, the attacker does not control the error. If the attacker controls the error locations, then it is possible to obtain x in $\lceil n/\varepsilon \rceil$ queries. This can happen during a fault attack, akin to side-channel attacks. It should also be noted that some coordinates of biometric data may be non-variable and, as a consequence, an attacker cannot recover them. This partial recovery attack is, therefore, a privacy attack, and leads to an authentication attack if the number of variable coordinates is sufficiently large (at least $n - \varepsilon$ in the binary case).

Remark 2.4.4. *In the non-binary case, the value $x_i - y_i$ does not provide enough information. The exact value of x_i can be determined in two cases. First, if $x_i - y_i = -q + 1$, then $x_i = 0$. Second, if $x_i - y_i = 2(q - 1)$, then $x_i = q - 1$. For all other cases, there is an ambiguity regarding the value of x_i as y_i is unknown. However, by knowing the distribution of x_i and y_i , repeating observations yields a statistical attack.*

Attacks for each type of leakage along with their complexities are summarized in Table 2.7.

| Distance-to-Threshold comparison | Leakage | Complexity type | Complexity in Big-O | Theorem |
|----------------------------------|-------------------------------------|-------------------------|---------------------|---------|
| Below | Distance | Exponential | $q^{n-\varepsilon}$ | 2.4.1 |
| | Positions | Exponential | $q^{n-\varepsilon}$ | 2.4.2 |
| | Positions and values | Exponential | $q^{n-\varepsilon}$ | 2.4.3 |
| | Positions and values (accumulation) | Linearithmic/Polynomial | $n^\alpha \log n$ | 2.4.7 |
| Both | Minimal | Exponential | $q^{n-\varepsilon}$ | [83] |
| | Distance | Linear | $(n-\varepsilon)q$ | 2.4.4 |
| | Positions | Constant | q | 2.4.5 |
| | Positions and values | Constant | 1 | 2.4.6 |

Table 2.8 – Summary of all leakage exploits and their complexities for the weaker attack model with α such that the occurrence of the rarest error is $n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$. The Distance-to-Threshold comparison determines if the leak occurs when $d(x, y) \leq \varepsilon$ (below) or when there is no distance requirement between x and y (both). For all the complexities, x and y are in \mathbb{Z}_q^n with $q \geq 2$ except for the minimal leakage where x and y are in \mathbb{Z}_2^n . The provided complexities represent worst-case scenarios, except for the accumulation attack where the result is the expectation.

¹Note that the Big-O complexity of the optimal exhaustive search strategy, in the worst-case, is the same as the naive strategy as the minimum of $h(\cdot)$ is 0.

2.4.6 Weaker Attack Model: Compromise a point in the ball

In the context of fuzzy data, some argue that identifying a nearby point is sufficient, as retrieving the exact data has no particular interest, given that this is already a fuzzy reading from the source. Consequently, within this particular framework, the center search attack presented by Pagnin *et al.* becomes irrelevant. In this section, we provide a comprehensive overview of the complexities associated with the discussed attacks straightforwardly derived from the previous sections in Table 2.8.

2.5 Untargeted Attack

The logical expansion of the study of targeted attacks is the analysis of untargeted attacks, as exemplified by password analysis in cryptography. The current configuration is the same as the one described in Section 2.4, with the exception that the adversary's objective is to impersonate a user of the database in the minimum number of rounds.

In the following, $\mathcal{B} = (v_1, \dots, v_N)$ denotes a template database and each $v \in \mathcal{B}$ is assumed to be uniformly drawn in \mathbb{Z}_2^n and independently from each other: $v \stackrel{\text{ind}}{\sim} \text{Unif}(\mathbb{Z}_2^n)$.

2.5.1 Metric Space Based Bounds

Below are presented some untargeted attacks to find near-collisions with hidden templates of a secret biometric database. We examine two attack scenarios, estimating the bounds for their respective run-time complexities. In the first scenario, an outsider attacker submits guesses to the system until one of them is accepted. In the multi-outsider scenario, several attackers launch an attack in parallel and try to impersonate any user. They apply independently of the operation mode (identification or verification). However, unlike identification, authentication requires a set of identifiers (logins). In this case, the attacker needs to test a guessed template for each of the identifiers, hence adding a factor of N in the estimated bounds. In fact, in the authentication case, the attacker performs one test per user instead of one in the identification case, hence the N factor.

Naive and Adaptive Attack Models: The Study of the Outsider

The attacker \mathcal{A} is an outsider of the system, *i.e.*, she is not enrolled in the system, and she seeks to perform an untargeted attack by impersonating any of the N users in the database. For the considered attack, the database is not leaked in cleartext, so the templates remain secret. It is assumed that \mathcal{A} generates her database composed of several templates $t_1^a, t_2^a, \dots, t_k^a$ using the transformation until one of them is accepted by the system. To do so, she generates templates randomly and applies the transformation to these templates. In the following, we calculate the probability of a strong near-collision for the generated template as well as the number of trials of a naive attacker.

Naive Attacker: We denote by E_1^N the event "the template of the outsider matches with at least one of the N templates of the database". When no constraints are imposed on enrollment templates, E_1^N is seen as a union of independent events E_1^1 . Consider that the attacker repeats this attack with each new generation of a template until achieving success. According to the geometric distribution, m_{out} the necessary number of templates to generate so that a success occurs with more than a chance of 50%, corresponds to the median number of trials to succeed, *i.e.*, about $-\log(2)/\log(1-p)$ where p is the probability of success with a trial.

Theorem 2.5.1. *Let N , n and ε be fixed parameters with $\varepsilon/n \leq 1/2$ and $N < 2^{n(1-h(\varepsilon/n))}$. The median number m_{out} of trials for the attacker to successfully impersonate a user is*

$$\Omega\left(2^{n(1-h(\varepsilon/n))-\log_2 N}\right) \text{ and } O\left(2^{n(1-h(\varepsilon/n))+\frac{1}{2}\log_2 \varepsilon(1-\frac{\varepsilon}{n})-\log_2 N}\right)$$

where $h(\cdot)$ is the binary entropy function.

Proof. The probability of success with a single trial is given by

$$p = 1 - (1 - V_\varepsilon)^N.$$

Since $-x - x^2 \leq \log(1-x) \leq -x$ for $0 \leq x \leq \frac{1}{2}$ and if $N < 2^{n(1-h(\varepsilon/n))}$, then

$$-N(V_\varepsilon - V_\varepsilon^2) \leq N \log(1 - V_\varepsilon) \leq -NV_\varepsilon$$

Next, since $V_\varepsilon^2 < V_\varepsilon$ and

$$\sum_{k=0}^{\varepsilon} \binom{n}{k} \leq 2^{nh(\varepsilon/n)}$$

for $\varepsilon/n < 1/2$ (see [119]), the number m_{out} of trials is lower bounded as follows:

$$m_{\text{out}} \geq \frac{\log 2}{N(V_\varepsilon + V_\varepsilon^2)} \geq \frac{\log 2}{2NV_\varepsilon} \geq \frac{\log 2}{2N} 2^{n(1-h(\varepsilon/n))} \quad (2.6)$$

For the upper bound, notice that for $\varepsilon/n < 1/2$, we have

$$\sum_{k=0}^{\varepsilon} \binom{n}{k} \geq \frac{2^{nh(\varepsilon/n)}}{\sqrt{8\varepsilon(1-\varepsilon/n)}}$$

(see [119]). Hence,

$$m_{\text{out}} \leq \frac{\log 2}{NV_\varepsilon} \leq \frac{\log 2}{N} 2^{n(1-h(\varepsilon/n))} \times \sqrt{8\varepsilon(1-\varepsilon/n)} \quad (2.7)$$

■

Case of Different Enrolled Templates: In this case, a verification is performed on the enrolled templates to ensure that the database is only comprised of distinct templates. In other words, it is necessary to consider that templates are dependent, which is an essential change compared to the context of Theorem 2.5.1. It is then worth noting that E_1^N cannot be seen as a union of independent events E_1^1 and that an exact measure of E_1^N involves the cardinalities of multiple intersections of Hamming balls. Therefore, the following result is a declination of Theorem 2.5.1 in this specific context.

Corollary 2.5.1. *Considering a similar setting of Theorem 2.5.1 but with distinct templates in the database, then the median number m_{out} of trials for the attacker to successfully impersonate a user is*

$$\Omega\left(2^{n(1-h(\varepsilon/n))-\log_2(N)-\log_2\left(1+6\frac{N-1}{2^{n+1}}\right)}\right)$$

and

$$O\left(2^{n(1-h(\varepsilon/n))+\frac{1}{2}\log_2\varepsilon\left(1-\frac{\varepsilon}{n}\right)-\log_2 N-\log_2\left(1+\frac{N-1}{2^{n+1}}\right)}\right).$$

Proof. Similarly to the proof of Theorem 2.5.1, the probability of success of a given trial is

$$p = \mathbb{P}\left(t \in \bigcup_{k=1}^N B_\varepsilon(v_k)\right) = 1 - \mathbb{P}\left(t \in \bigcap_{k=1}^N \overline{B_\varepsilon(v_k)}\right)$$

where t is the generated template of the attacker, and v_k is the k -th enrolled template in the template database. As templates v_1, \dots, v_N are not independent, an alternative formulation can be expressed with conditional probabilities. Each conditional probability corresponds to the event that the $(\ell + 1)$ th enrolled template is sampled without replacement and that it does not be matched with t :

$$\begin{aligned} \mathbb{P}\left(t \in \bigcap_{k=1}^N \overline{B_\varepsilon(v_k)}\right) &= \prod_{\ell=0}^{N-1} \mathbb{P}\left(t \in \overline{B_\varepsilon(v_{\ell+1})} \mid t \in \bigcap_{k=1}^{\ell} \overline{B_\varepsilon(v_k)}\right) \\ &= \prod_{\ell=0}^{N-1} \frac{2^n - \ell - |B_\varepsilon|}{2^n - \ell} = \prod_{\ell=0}^{N-1} \left(1 - \frac{2^n}{2^n - \ell} V_\varepsilon\right) \end{aligned}$$

Next, since $-x - x^2 \leq \log(1-x) \leq -x$ for $0 \leq x \leq \frac{1}{2}$ and if $N < 2^{n(1-h(\varepsilon/n))-1}$, then

$$-V_\varepsilon S_1 - V_\varepsilon^2 S_2 \leq \log(1-p) \leq -V_\varepsilon S_1$$

where $S_1 = \sum_{\ell=0}^{N-1} \frac{2^n}{2^n - \ell}$ and $S_2 = \sum_{\ell=0}^{N-1} \left(\frac{2^n}{2^n - \ell}\right)^2$. Moreover, for what follows, S_1 and S_2 can be bounded as

$$S_1 \geq N \left(1 + \frac{N-1}{2^{n+1}}\right) \quad \text{and} \quad S_2 \leq N \left(1 + 6\frac{N-1}{2^{n+1}}\right)$$

since $(1-x)^{-1} \geq 1+x$ and $(1-x)^{-2} \leq 1+6x$, for $0 \leq x \leq \frac{1}{2}$. Lastly, notice that $S_2 > S_1$ and then the results follow. ■

Some care should be taken for the choice of ε , since a high value for ε dramatically reduces the

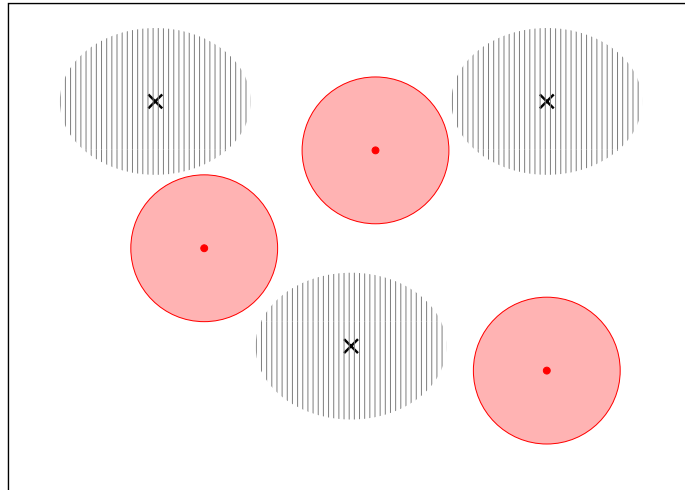


Figure 2.9 – Representation of a κ -adaptive seeking a template of a user. The cross symbols indicate unsuccessful attempts, while the hashed areas represent the knowledge gained about the absence of templates in those regions.

number of trials of the attacker, despite a large n . Theorem 2.5.1 makes clear the link between security parameters, hence allowing a secure choice of the parameters.

κ -adaptive Attacks: To study a case close to what might be a smart attacker, we consider that after a trial the attacker can infer that some of the remaining templates are not a right fit. Next candidate guesses that are FMR away from tried templates are better choices than those near the center. Such an inference should vary depending on the number of trials made and potentially extra information and should lead to different amounts of inferred templates. For the sake of simplicity, we investigate below an attacker model for which the attacker infers in average κ unsuccessful templates. Figure 2.9 gives a representation of an attack performed by the κ -adaptive attacker.

Definition 2.5.1 (κ -adaptive attacker). *An attacker is κ -adaptive if for each of its trials to impersonate a user template she is able to identify κ non-hit templates.*

As an example, a 0-adaptive attacker tries to impersonate a user template by sampling into \mathbb{Z}_2^n with replacement, and sampling without replacement for a 1-adaptive attacker. Proposition 2.5.1 states that under reasonable conditions for the parameters, the number of tries required for a κ -adaptive attacker to succeed is equivalent to the required number of tries for a 0-adaptive attacker. Let $A(\kappa)$ denote the number of trials of a κ -adaptive attacker to succeed.

Proposition 2.5.1. *If p is negligible and κ negligible compared to $\sqrt{2^n}$ then, for a given number of trials $a \leq \sqrt{2^n}$, the probability that, among an amount of "a" trials, at least one successful trial of a κ -adaptive attacker is asymptotically equivalent to at least one successful trial of a 0-adaptive attacker:*

$$\mathbb{P}(A(0) \leq a) \sim \mathbb{P}(A(\kappa) \leq a).$$

Proof. Let p be a negligible probability, and notice that it is the case when the number of user templates is negligible compared to 2^n . Then, $\kappa = g(n) \in o(2^{n/2})$ and $a \in \{1, \dots, 2^{n/2}\}$. The probability of $\mathbb{P}(A(0) = a)$ is given by a geometric law of probability p . Thus, $\mathbb{P}(A(0) \leq a) = p \sum_{i=1}^a (1-p)^{i-1}$ and

by using Lemma 2.5.1, we have:

$$\mathbb{P}(A(\kappa) \leq a) = \sum_{i=1}^a \frac{p}{1 - (i-1)\frac{\kappa}{2^n}} \times \prod_{j=1}^{i+2} \frac{1 - p - j\frac{\kappa}{2^n}}{1 - j\frac{\kappa}{2^n}}.$$

When n tends to infinity, according to the assumptions, the result follows since

$$\frac{\mathbb{P}(A(g(n)) \leq a)}{p} \xrightarrow{n \rightarrow +\infty} a \quad \text{and} \quad \frac{\mathbb{P}(A(0) \leq a)}{p} \xrightarrow{n \rightarrow +\infty} a.$$

Then, as the sums are finite, the result follows. \blacksquare

Proposition 2.5.2 states that under reasonable conditions for the parameters, the probability that a κ -adaptive attacker succeeds at the a^{th} trial is equivalent to the probability that a 0-adaptive attacker succeeds at the same trial.

Proposition 2.5.2. *If p is negligible and κ negligible compared to $\sqrt{2^n}$ then, for a given number of trials $a \leq \sqrt{2^n}$, the probability of a successful trial of a κ -adaptive attacker is asymptotically equivalent to the probability of successful trial of a 0-adaptive attacker:*

$$\mathbb{P}(A(0) = a) \sim \mathbb{P}(A(\kappa) = a).$$

In order to provide a proof, we need two intermediate results. Lemma 2.5.1 gives the probability that the first success of a κ -adaptive attacker is the a^{th} trial.

Lemma 2.5.1. *The probability that the first success of a κ -adaptive attacker is the a^{th} trial is given by*

$$\mathbb{P}(A(\kappa) = a) = \frac{p2^n}{2^n - \kappa(a-1)} \frac{\binom{\frac{2^n(1-p)}{\kappa}}{a-1}}{\binom{\frac{2^n}{\kappa}}{a-1}}$$

for $a \in \{1, \dots, \lceil \frac{2^n(1-p)}{\kappa} \rceil + 1\}$ and for $\kappa > 1$.

Proof. If none of the $(a-1)^{\text{th}}$ trials lead to success, then the probability that the a^{th} trial fails, is

$$\begin{aligned} \mathbb{P}(A(\kappa) > a \mid A(\kappa) > a-1) &= \frac{2^n(1-p) - (a-1)\kappa}{2^n - (a-1)\kappa} \\ &= \frac{\frac{2^n(1-p)}{\kappa} - (a-1)}{\frac{2^n}{\kappa} - (a-1)} \end{aligned}$$

then

$$\begin{aligned} \mathbb{P}(A(\kappa) = a) &= \left(1 - \mathbb{P}(A(\kappa) > a \mid A(\kappa) > a-1)\right) \times \prod_{j=1}^{a-1} \mathbb{P}(A(\kappa) > j \mid A(\kappa) > j-1) \\ &= \frac{p2^n}{2^n - \kappa(a-1)} \times \frac{\frac{2^n(1-p)}{\kappa}}{\frac{2^n}{\kappa}} \times \dots \times \frac{\frac{2^n(1-p)}{\kappa} - (a-2)}{\frac{2^n}{\kappa} - (a-2)} \end{aligned}$$

and the result follows since parts of the products above correspond to gamma function rates, which are related to binomial coefficients. \blacksquare

As we have the probability that the first success of a κ -adaptive attacker is the a^{th} trial and the probability that the first success of a 0-adaptive attacker is the a^{th} trial, Lemma 2.5.2 investigate the ratio between those two to compare them.

Lemma 2.5.2. *The rate of probability of success at a given trial between an 0-adaptive attacker and a κ -adaptive attacker is*

$$\frac{\mathbb{P}(A(0) = a)}{\mathbb{P}(A(\kappa) = a)} = \prod_{j=1}^{a-1} \frac{1 - j \frac{\kappa}{2^n}}{1 - (j-1) \frac{\kappa}{2^n(1-p)}}$$

for $a \in \{1, \dots, \lceil \frac{2^n(1-p)}{\kappa} \rceil + 1\}$.

Proof. Concerning the 0-adaptive attacker, which corresponds to a sampler with replacement, the probability of a first success is a geometric distribution:

$$\mathbb{P}(A(0) = a) = p(1-p)^{a-1}.$$

According to Lemma 2.5.1, we have

$$\begin{aligned} \frac{\mathbb{P}(A(0) = a)}{\mathbb{P}(A(\kappa) = a)} &= (1-p)^{a-1} \frac{2^n - \kappa(a-1)}{2^n} \frac{\binom{\frac{2^n}{\kappa}}{a-1}}{\binom{\frac{2^n(1-p)}{\kappa}}{a-1}} \\ &= (1-p)^{a-1} \left(1 - (a-1) \frac{\kappa}{2^n}\right) \times \frac{\prod_{j=1}^{a-2} (1 - j \frac{\kappa}{2^n})}{\prod_{j=0}^{a-2} (1 - p - j \frac{\kappa}{2^n})} \\ &= \prod_{j=1}^{a-1} \frac{1 - j \frac{\kappa}{2^n}}{1 - (j-1) \frac{\kappa}{2^n(1-p)}}. \end{aligned}$$

■

With these two intermediate results, we can now prove Proposition 2.5.2.

Proof. Let p be a negligible probability, and notice that it is the case when the number of user templates is negligible compared to 2^n . Then, $\kappa = g(n) \in o(2^{n/2})$ and $a \in \{1, \dots, 2^{n/2}\}$. Using Lemma 2.5.2, we have

$$\frac{\mathbb{P}(A(0) = a)}{\mathbb{P}(A(\kappa) = a)} = \prod_{j=1}^{a-1} \frac{1 - j \frac{\kappa}{2^n}}{1 - (j-1) \frac{\kappa}{2^n(1-p)}}.$$

When n tends to infinity, according to the above assumptions, $\forall j \in \{1, \dots, a-1\}$

$$1 - j \frac{g(n)}{2^n} \xrightarrow{n \rightarrow +\infty} 1 \quad \text{and} \quad 1 - (j-1) \frac{g(n)}{2^n(1-p)} \xrightarrow{n \rightarrow +\infty} 1.$$

As all the terms of the product are positive, the result follows. ■

From Proposition 2.5.1, we show that the median number of trials is equivalent for both attacker models.

Theorem 2.5.2. *Let m_κ the median number of trials of a κ -adaptive attacker to obtain a first success. If p is negligible, κ negligible compared to $\sqrt{2^n}$ and $m_\kappa \leq \sqrt{2^n}$, then, $m_0 \sim m_\kappa$.*

| | | | | | | |
|--|--------|-----|-----|-----|--------|--------|
| n | 128 | 256 | 128 | 128 | 128 | 128 |
| ε | 30 | | 10 | 20 | 30 | 30 |
| $N (\log_{10})$ | 6 | | 6 | | 5 | 7 |
| $a (\log_{10})$ | 3 | | 3 | | 3 | 2 4 |
| $\kappa (\log_2)$ | 47 | | 47 | | 47 | 47 |
| $\frac{\mathbb{P}(A_h(\kappa) \leq a)}{\mathbb{P}(A_h(0) \leq a)}$ | 0.9981 | 1 | 1 | | 0.9998 | 0.9808 |
| | | | | | | 20 100 |
| | | | | | | 0.9981 |
| | | | | | | 0.9981 |

Table 2.9 – Ratio between a κ -adaptive attacker and a 0-adaptive attacker in function of n , ε , N , a and κ for the lower bound.

Proof. According to Proposition 2.5.1,

$$\frac{\mathbb{P}(A(0) \leq m_\kappa)}{\mathbb{P}(A(\kappa) \leq m_\kappa)} \xrightarrow{n \rightarrow +\infty} 1$$

and since $\mathbb{P}(A(\kappa) \leq m_\kappa) = 1/2$, we derive that m_0 tends to m_κ as n increases, if $m_\kappa \leq \sqrt{2^n}$. ■

A κ -adaptive attacker is under realistic constraints equivalent to a 0-adaptive attacker. Thus, for the sake of simplicity, considering Proposition 2.5.1, Theorem 2.5.2 and Proposition 2.5.2, in the remainder of the section, we only derive theoretical results for a 0-adaptive attacker model.

Remark 2.5.1. *In practice, an attacker generates a template from the set \mathbb{Z}_2^n deprived of the previously generated templates. When the number of trials (i.e., rounds) until the first success is low, the proposed simplified experiment above has a negligible bias. The larger is N , the lesser is the number of trials until a first success. In a biometric database, the number N of templates can be assumed large enough so that the cardinal of \mathbb{Z}_2^n overwhelms the number of trials.*

Numerical Results

To support Proposition 2.5.1 considering finite parameters, Table 2.9 investigates several settings by computing the ratio

$$\frac{\mathbb{P}(A(\kappa) \leq a)}{\mathbb{P}(A(0) \leq a)}$$

as well as $\mathbb{P}(A(\kappa) \leq a)$. In the following, the union size is maximized (i.e., all the balls are disjoint). Hence, a case in favor of the attacker and against our proposal is investigated. One of our realistic settings is defined as ($n \geq 128, \varepsilon \leq 0.25n, N = 10^6, \kappa \leq |B_\varepsilon|, a = 10^3$) where a and κ are respectively the number of trials and the extra information of a κ -adaptive attacker. In this case, we observe a ratio very close to 1. It is interesting to note that even with a large increase in the information given to the κ -adaptive attacker, her probability of success does not increase significantly. Hence, as shown by Proposition 2.5.1, Proposition 2.5.2, Theorem 2.5.2 and Table 2.9, a 0-adaptive attacker performs as well as a κ -adaptive attacker for reasonable parameters.

Multiple Attackers

In order to provide a comprehensive analysis of the various scenarios that may arise in this context, we consider the case where multiple independent attackers target the same system. The success of the attackers is then evaluated by determining whether at least one of them succeeds in

| n | 128 | | | | | | | | | | | | 256 | | | | | | | | | | | | | | | |
|-----------------------------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $N (\log_{10})$ | 4 | | | | 6 | | | | 8 | | | | 4 | | | | 6 | | | | 8 | | | | | | | |
| ε | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 | 12 | 19 | 25 | 51 |
| $h(\varepsilon/n)$ | 0.4 | 0.6 | 0.7 | 1.0 | 0.4 | 0.6 | 0.7 | 1.0 | 0.4 | 0.6 | 0.7 | 1.0 | 0.3 | 0.4 | 0.5 | 0.7 | 0.3 | 0.4 | 0.5 | 0.7 | 0.3 | 0.4 | 0.5 | 0.7 | 0.3 | 0.4 | 0.5 | 0.7 |
| Lower bound Outsider (\log_2) | 62 | 45 | 33 | 8 | 56 | 38 | 27 | 1 | 49 | 32 | 20 | 0 | 178 | 153 | 135 | 78 | 171 | 147 | 128 | 72 | 165 | 140 | 122 | 65 | | | | |
| Upper bound Outsider (\log_2) | 64 | 47 | 36 | 10 | 57 | 40 | 29 | 4 | 51 | 34 | 22 | 0 | 180 | 155 | 137 | 81 | 173 | 149 | 131 | 74 | 166 | 142 | 124 | 68 | | | | |

Table 2.10 – Bounds for the number of operations for an outsider, in function of n , N and ε .

matching. It should be noted that, as these are independent attackers, they do not communicate with each other. Furthermore, we consider attackers to be 0-adaptive, so the fact that different attackers each try the same template corresponds to the case of the naive attacker who does not adapt his strategy. In other words, having several attackers in this case is equivalent to having a single attacker who tries several times per round. In the following analysis, we assume that the number of attackers is proportional to the number of individuals in the database. It is reasonable to assume that the greater the number of individuals in the database, the greater the number of attackers who may be interested in compromising the database. Alternatively, this may be interpreted as an assumption that the total number of users may include a proportion of malicious users.

In what follows, α is used to denote the proportion of malicious users, and it is assumed that there are αN attackers in this scenario. Below, Theorem 2.5.3 provides an asymptotic result about m_α , the median number of rounds until at least one attacker successfully impersonates a user of the database.

Theorem 2.5.3. *Let n the template size and let N and ε be fixed parameters with $\varepsilon/n \leq 1/2$ and $N < 2^{n(1-h(\varepsilon/n))-1}$. If $\frac{-\log 2}{N^2 \log(1-V_\varepsilon)} \geq \alpha > 0$ is the portion of malicious users, then the median number m_α of rounds until at least one attacker succeeds is*

$$\Omega(2^{n(1-h(\varepsilon/n))-2\log_2(N)-\log_2(\alpha)}) \quad \text{and} \quad O(2^{n(1-h(\varepsilon/n))+\frac{1}{2}\log_2 \varepsilon(1-\frac{\varepsilon}{n})-2\log_2(N)-\log_2(\alpha)}).$$

Proof. At each round, the probability of a success of the i^{th} attacker is $p_i = 1 - (1 - V_\varepsilon)^N$, as described in the proof of Theorem 2.5.1. Since the attackers are independent, then the probability of success at each round is $p = 1 - \prod_{i=1}^{\alpha N} (1 - V_\varepsilon)^N$, then $p = 1 - (1 - V_\varepsilon)^{\alpha N^2}$. Next, according to the same procedure as in the proof of Theorem 2.5.1, we have

$$\frac{\log 2}{\alpha N^2} 2^{n(1-h(\varepsilon/n))} \leq m_\alpha \leq \frac{\log 2}{\alpha N^2} 2^{n(1-h(\varepsilon/n))} \times \sqrt{8\varepsilon(1-\frac{\varepsilon}{n})}$$

and the result follows. ■

2.5.2 Biometric security: Assess the Security Score Against the Outsider

In this section, Proposition 2.5.1 is evaluated on reasonable scenarios, providing the number of guesses needed for an outsider attacker to impersonate a user along with a new security metric for evaluating the resilience of a database with respect to this scenario.

The analyses focus on a database \mathcal{B} of uniformly distributed template in \mathbb{Z}_2^n , *i.e.*, when the biometric transformation performs like a perfect random function. This enables to provide an upper bound on run-time complexities. Although the assumption of a uniform distribution yields

an overestimated upper bound, it is helpful for securely parametrizing the transformation scheme. However, concerning the lower bound, it is above reality if the distribution is non-uniform, as it is the case for deterministic biometric transformation.

We consider the resilience to attacks in the outsider scenarios as an indicator of the security of biometric databases. In Table 2.10, some numeric values for n , N , and ε give a poor resistance to those attacks. The score S_1 is introduced accordingly. By denoting p_1 the lower bounds for the number of rounds of an outsider, the corresponding score is $S_1 = \log_2(p_1/2^{128})$. The database is resilient to these attacks if the scores are greater than or equal to 0. According to Table 2.10, the triplet ($n = 256$, $N = 10^4$, $\varepsilon = 19$) yields the scores $S_1 = 25$. However, the triplet ($n = 128$, $N = 10^4$, $\varepsilon = 12$) yields the scores $S_1 = -66$. Usually, biometric recognition systems are parameterized by adjusting a threshold using a training dataset in experimental evaluations to find the Equal Error Rate (EER). However, the obtained threshold could be too large, whence not providing the expected level of security based on our analyses. Therefore, a trade-off has to be found between the False Non-Matching Rate (FNMR, with respect to the EER) and the above security scores.

2.6 Conclusion and Future Work

In this chapter, we investigate the impact of near-collisions and master-template events on the security of biometric systems. As these events are complex to consider, we approximate their probabilities for a database of known and unknown templates. Additionally, we introduce several attacker scenarios to provide a number of interpretable results that could inform the overall security of a biometric system. The results demonstrate the significance of selecting appropriate parameters for a database. They offer a methodology for meticulously selecting these parameters. Furthermore, the study introduces a novel adaptive security metric, based on these attacks, and examines the probability of a weak near-collision and of a master template occurrence. For future research, all the results provided may be generalized fairly easily to the case $q > 2$ by using the cardinal of a ball in \mathbb{Z}_q^n given by

$$B_q = \sum_{i=0}^{\varepsilon} \binom{n}{i} (q-1)^i$$

and the bounds given by

$$q^{nh_q(\varepsilon/n)-o(n)} \leq B_q \leq q^{nh_q(\varepsilon/n)}$$

with h_q the q -ary entropy function [72]. In addition, to expand the scope of the current findings we may consider other distances. Another direction for further investigation is to broaden these results to the case where the templates are not uniformly distributed. In the non-uniform case, specifically within the binary space \mathbb{Z}_2^n , we may use a multivariate Bernoulli model to capture the fact that certain bits may occur more frequently than others, and highlight groups of bits that may be statistically dependant [86]. This tool could be employed for both theoretical and practical analysis using real template databases. A practical analysis could reveal blocks of correlated bits, which would have an impact on the security of the system.

Subsequently, we have examined the impact of targeted and untargeted attacks on the security of the systems. Our investigation into the information leakage of a biometric system using privacy-preserving distance has revealed critical security vulnerabilities that arise under various scenarios. By evaluating the impact of different types of leakage, including distance, error position, and error value, we have examined the potential risks posed to data privacy and security in practical

applications. It is important to note that leakage below the threshold does not significantly degrade the security of the system. In contrast, leakage above and below the threshold significantly decreases the security. The accumulation attack that we introduce and investigate is based on the assumption that errors are uniformly distributed throughout each authentication session. The result of the accumulation attack may be further refined by considering a variable number of coupons, randomly drawn between 0 and ε in each round, while acknowledging the actual distribution of the errors. Further research should focus on refining the accumulation attack as previously suggested. Given the inherent complexity of this probabilistic problem, it is imperative to undertake preliminary empirical analysis through simulation. The aforementioned simulations would yield a relatively precise estimation of the number of rounds if a sufficient number of simulations are performed. One drawback of this methodology is the necessity to ascertain the distribution of errors, which may vary considerably depending on the system used. This could result in a non-generic outcome. More precisely, the number of rounds for a given system would be determined, it is not possible to generalize this to all existing systems using the same template size and threshold as the error distribution may harshly change. Additionally, it would be beneficial to identify instances of the accumulation attack being applied to systems in the literature and to make it more generic. The estimation of the proportion of attackers on a database would also enable the refinement of results in the case of exhaustive searches with multiple attackers.

To achieve a comprehensive understanding of the security implications of these attacks, it is necessary to develop additional metrics that can be used to effectively compare the security of different systems and determine their security in an interpretable way. Therefore, the security and performance of biometric systems would be multifactorial and not merely reflected in terms of false match rate FMR or FAR.

Design of a Remote Secure Sketch Resistant to Offline Attacks

Outline of the current chapter

| | |
|--|-----------|
| 3.1 Introduction | 64 |
| 3.2 Generic Construction of Secure Sketches from Groups | 67 |
| 3.2.1 Mathematical Construction: Groups with Unique t -factorization | 67 |
| 3.2.2 Computational Secure Sketch Construction | 70 |
| 3.2.3 Security of SPS | 72 |
| 3.2.4 On the Hardness of and Relationships between DSPP and CSPP | 75 |
| 3.2.5 Robust SPS in the Random Oracle Model | 84 |
| 3.2.6 Concrete Instantiations of SPS | 85 |
| 3.2.7 Details on VBB and NTT | 85 |
| 3.2.8 Robust Fuzzy Extractor based on Robust SPS | 88 |
| 3.3 Biometric Authentication Protocol from Groups | 89 |
| 3.4 Remote Secure Sketch from Group | 92 |
| 3.4.1 Requirements for Building a RSS from Subset Product Problems | 93 |
| 3.4.2 Mechanics of our Remote Subset Product Sketch | 93 |
| 3.4.3 Authentication and Secrecy of Fresh Readings | 96 |
| 3.4.4 Secrecy of the Template Sketch | 97 |
| 3.5 Secure Remote Fuzzy Extractors from Remote Secure Sketches | 98 |
| 3.6 Conclusion and Future Work | 99 |

Abstract of the current chapter:

In this chapter, we focus on how to securely use biometric data. Specifically, we aim to avoid or limit offline exhaustive search attacks, as these are more devastating than online attacks. Online attacks are generally detected or mitigated through counter-measures such as increasing wait times between attempts. Our objective is to construct a cryptographic protocol that allows user authentication based solely on biometric data while protecting against offline attacks. We proceed in two steps, on the first hand, we develop a mathematical structure, and on the other hand, we design two distinct protocols based on this structure.

First, we developed a mathematical structure to evaluate the distance between two vectors in an homomorphic manner. This structure also allows for the correction of the provided data to match the reference data if the distance is below a given threshold. The main advantages of this construction lie in that recovering the reference data and linking two enrollments are computationally hard.

Then, with this construction, we developed two protocols serving distinct purposes. The first protocol fully addresses our original problem by using a *zero-knowledge* proof mechanism to ensure that a computationally unbounded malicious client cannot obtain useful information during the authentication. In addition, we propose a second protocol that ensures both authentication and reconstruction of the original secret. Note that this protocol is only suitable if the adversary is assumed to be computationally bounded, which implies a lower level of security. However, both versions ensure a high level of security against offline attacks, as these attacks are either not possible or have an exponential cost.

3.1 Introduction

Context: A typical online biometric authentication protocol comprises two phases, with the client and server engaged in a series of interactions. The enrollment and verification phases are typical of an online biometric authentication protocol. During the enrollment phase, biometric samples are collected from the user, which are then processed to derive a biometric template using a feature extraction algorithm. The server stores the biometric template (or some information derived from the template, *e.g.*, a cryptographic key) together with the user's identifier (ID). In the verification phase, the user regenerates their biometric template and uses it in the protocol to prove the authentic ownership of their biometric data (enrolled under their ID) against the server.

The processing, transmission, and storage of information derived from users' biometric data, also known as the biometric template, is a fundamental aspect of online biometric applications. Biometric templates serve as the primary reference data for recognizing individuals uniquely in applications and are considered part of personally identifiable information. Consequently, the protection of individual's biometric information and their privacy is of paramount importance in biometric systems and applications. Some regulations, such as the General Data Protection Regulation (GDPR) in the EU and the California Consumer Privacy Act (CCPA) for California residents of the US, have been enacted to protect individuals' data and privacy. The inherently noisy nature of biometrics, coupled with the fact that the biometric characteristics of an individual are not easily renewable, makes the design of secure biometric-based authentication protocols a more challenging endeavor than the design of token or password-based protocols. Research

and standardization efforts [23, 1, 2] have identified several requirements for securing biometric information and templates, including *renewability*, *irreversibility*, *unlinkability*, *indistinguishability*, and *reusability*. Informally, *renewability* is the ability to create (randomized) biometric templates from the same biometric data. *Irreversibility* implies that it is infeasible to recover the plain biometric data from its protected template. *Unlinkability*, *indistinguishability*, and *reusability* are closely related, and they require that an adversary, who observes a user's protected biometric templates enrolled at different servers, cannot yield a significant advantage towards degrading that user's security or privacy, such as cross-matching the individual's records or recovering their biometric data by reversing their biometric templates (See Section 1.3.1 for more detail).

Secure sketch schemes represent a primary cryptographic primitive utilized to protect biometric templates with formal security guarantees. Informally, given a biometric input vector x , a sketch S is derived through a randomized process. In this context, the sketch S should be irreversible, yet it should allow the recovery of x in the presence of another biometric input vector y that is close to x . In this chapter, we show a generic method for constructing a family of secure sketches from a specific family of groups. We refer to these secure sketches as the Subset Product Sketch (SPS) because the construction relies on multiplying group elements from a particular subset. Secure sketches are designed to allow the recovery of originally enrolled data from inputs that may vary slightly over time. This capability is of significant importance in applications where data consistency cannot be guaranteed due to natural variations, such as in biometric systems and hardware security. Traditionally, secure sketches are constructed using error-correcting codes [122, 120, 109] to effectively handle these variations. A secure sketch can be viewed as a form of error-correcting code. Error-correcting codes [145, 98, 146, 147] aim to recover the original data from a noisy version as secure sketches. Thus, both concepts share a common philosophy: the correction of errors. However, secure sketches aim to preserve the privacy of the enrolled vector, which is not necessary the case of an error-correcting code. Principles of information theory ensure the security of these sketches by managing the tradeoff between data recoverability and confidentiality. More specifically, the information revealed via publishing the sketch value of a secret input is required to be bounded for the (*average*) *min-entropy* notion [96]. A relaxation has been made in [40] by switching from *min-entropy* to *Hill-entropy* [96], which mainly says that the distribution of secret inputs given their sketch values have entropy at least k if that distribution is computationally indistinguishable from a distribution (conditioned on their sketch values) with entropy at least k with respect to the *min-entropy*. We avoid associating entropy-based security notions to sketch schemes and base our security on the hardness of decisional and computational problems.

Detailed Contributions: In this chapter, we demonstrate how to construct a "Subset-Product Sketch" (SPS) and how to use it for implementing a remote secure sketch. SPS functions as a secure sketch, meaning it allows the recovery of a stored secret by providing a similar secret within a certain distance and threshold, using helping data. This helping data, generated from the secret, must not reveal any information about the secret itself. Additionally, SPS must satisfy the properties of being reusable and irreversible. Reusable means that if multiple helping data are derived from the same secret, collecting these data does not provide more information about the secret than a single one. Irreversible means that given the helping data, it is difficult or impossible to recover the original secret.

To achieve this, we first define a family of groups with a unique factorization property. This unique factorization property allows us to hide a secret in the exponent of a base and, under certain

conditions, to uniquely recover this exponent. More precisely, under the Gaussian Heuristic, this property enables the unique correction of errors within the verification threshold. Additionally, hiding the secret in the exponent ensures irreversibility under hard computational assumptions such as the Computational Subset Product Problem (CSPP). This construction allows for the random selection of groups and bases, ensuring reusability under decision problems such as the Decisional Subset Product Problem (DSPP). Thus, through this construction, SPS can be implemented and ensure all the aforementioned properties. We emphasize that our SPS can create sketches from integer vectors and tolerate a linear fraction of errors with respect to the L_1 distance (*i.e.*, for $x, y \in \mathbb{Z}^n$, $L_1(x, y) = \sum_{i=0}^n |x_i - y_i|$). This functionality surpasses some previous constructions, which are limited to binary vectors and Hamming distance [63, 62, 60, 33].

Finally, we demonstrate how to construct remote secure sketches and remote fuzzy extractors based on SPS. Our first protocol is a zero-knowledge protocol, utilizing the checkability of the factorization correctness. This construction ensures total resistance to offline attacks, as a malicious client does not obtain any targets for their attacks. In the second protocol, we allow the client to correct its data while ensuring resistance to offline attacks. In this case, since the client receives obfuscated helping data, the complexity of exhaustive search is guaranteed by the aforementioned computational assumptions and the difficulty of breaking the obfuscation.

A major part of those results are from Durbet *et al.* [12].

Related works: Two closely related constructions are the virtual black box (VBB) [49] and the noise tolerant template (NTT) [71] schemes. They consist of deterministic functions and do not yield a sketch scheme as defined. Also, their security has not been previously analyzed with respect to reusability. As we discuss in Section 3.2.6, they can be used to realize concrete instantiation of SPS, and our generic construction provides a unified understanding of these primitives and their security as secure sketches. Early examples of secure sketches mostly rely on error-correcting codes, where the noise tolerance is measured with respect to Hamming distance or set difference metric, and their error tolerance is bounded by the error-correcting capacity of the underlying code. Fuzzy commitment [136] and fuzzy vault [116] schemes are two well-known constructions, and for a more extensive treatment of sketches and extractors based on error-correcting codes, we refer to [109]. In secure sketches and their extension to fuzzy extractor schemes, adversaries can exploit (distinct but correlated) sketches of the same client over different servers and gain significant information about their secret input. This is also known as the reusability attack [122, 108, 63]. Apon *et al.* [63] show that reusable fuzzy extractors can be constructed based on learning with errors problem (LWE [107]). However, [63] can tolerate a sublinear fraction of errors (as opposed to linear). Furthermore, [63] requires that some universal public domain parameters be used across all service providers which may not be practical for implementing the scheme in real-life applications. Another reusable fuzzy extractor is constructed in [33], where the idea is to sample a sufficiently large number of sufficiently small subsets from noisy data so that samples from a relatively close data pair contain at least one identical pair that can be verified using *digital lockers*. Due to the communication and memory cost, this scheme and its variants are not yet considered to be practical [54, 22]. Another disadvantage of [33] is its low error tolerance rate $k/(n \log n)$, where k is the length of the subsequences. Other reusable fuzzy extractors have been proposed based on LWE and discrete logarithm problems [60, 62, 61]. As discussed in this chapter, our secure sketch construction can handle a linear fraction of errors with respect to the Hamming and L_1 distances, and satisfy reusability and irreversibility, where best-known attacks seem to have

exponential complexity in the length of input vectors.

3.2 Generic Construction of Secure Sketches from Groups

In this section, we describe how to construct sketch schemes from certain families of groups. First, we define a specific family of groups \mathbb{G} , which have the unique t -factorization property with respect to a basis \mathbb{B} . Then, we show that the pair (\mathbb{G}, \mathbb{B}) (in a generic sense) yields sketch schemes, which we refer to as the Subset Product Sketch (SPS). In the end, the security of our scheme is discussed.

Intuition of the Construction: The idea is to conceal a vector by using its elements as exponents in the product of integers modulo another integer. If the product exceeds the modulo, the structure of the product is lost, making factorization useless for recovering the vector. However, if a product with the same elements raised to a power close to the original is provided, the division of these two products would be smaller than the modulo. By factoring the quotient of the products, we can retrieve the coordinates and the difference values between the enrolled vector and the tested vector. This error vector is then used to correct the tested data. The key point is that if the two vectors are not sufficiently close, the quotient remains larger than the modulo, preventing the factorization from yielding any useful information.

3.2.1 Mathematical Construction: Groups with Unique t -factorization

Let \mathbb{G} be a finite multiplicative group and let \mathbb{B} be a finite subset of \mathbb{G} , where all pairwise distinct elements $u_i, u_j \in \mathbb{B}$ are ordered using the lexicographic ordering on the binary representation of group elements. More precisely, we have $u_i < u_j$ for $i < j$. We also define

$$\mathbb{G}_{\mathbb{B}, t} = \left\{ \prod_{i=1}^{|\mathbb{B}|} u_i^{\delta_i} : u_i \in \mathbb{B}, \delta_i \in \mathbb{Z}, \sum_{i=1}^{|\mathbb{B}|} |\delta_i| \leq t \right\}. \quad (3.1)$$

Definition 3.2.1 (unique t -factorization property). \mathbb{G} has a unique t -factorization property with respect to a (factor) basis \mathbb{B} , if for all $g \in \mathbb{G}_{\mathbb{B}, t}$, there is a unique ordered integer sequence $[\delta_i]_{i=1}^{|\mathbb{B}|}$ such that

$$g = \prod_{i=1}^{|\mathbb{B}|} u_i^{\delta_i} \text{ and } \sum_{i=1}^{|\mathbb{B}|} |\delta_i| \leq t. \quad (3.2)$$

For \mathbb{G} with a unique t -factorization property, we associate an algorithm *Factor* that takes as input $g \in \mathbb{G}_{\mathbb{B}, t}$ and outputs $[\delta_i]_{i=1}^{|\mathbb{B}|}$ satisfying (3.2).

Remark 3.2.1. Even though we do not require *Factor* to be a polynomial time algorithm in this section, its efficiency becomes important when we define our subset product sketch scheme SPS. Section 3.2.6 provides two instantiations of SPS based on previous work [49, 71], where the complexity of *Factor* is polynomial in $\log |\mathbb{G}|$.

Example 3.2.1.1. $\mathbb{G} = \mathbb{Z}_{31}^*$ has a unique 2-factorization property with respect to the (factor) basis $\mathbb{B} = \{2, 3\}$ because

$$\mathbb{G}_{\mathbb{B}, 2} = \{1, 2, 3, 4, 6, 7, 8, 9, 11, 16, 17, 21, 26\} \quad (3.3)$$

has 13 elements and that there are exactly 13 distinct ordered integer sequences $[\delta_i]_{i=1}^2$ with $\sum_{i=1}^2 |\delta_i| \leq 2$. However, $\mathbb{G} = \mathbb{Z}_{31}^*$ does not have a unique 2-factorization property with respect to the (factor) basis $\mathbb{B} = \{2, 3, 11\}$ because there are two distinct ordered integer sequences $[\delta_i]_{i=1}^3$ such that $\sum_{i=1}^3 |\delta_i| \leq 2$, and that yield the same element. Namely, for $[1, 0, 0]$ and $[0, 1, 1]$, we have $2 = 2^1 3^0 11^0 = 2^0 3^1 11^1$ in \mathbb{G} .

Example 3.2.1.1 shows that groups with unique t -factorization property exist. Next, we show in Theorem 3.2.1 and Corollary 3.2.1 that for sufficiently large prime order groups with t -factorization property, the uniqueness property follows under the Gaussian heuristic [110].

Heuristic 3.2.1.1 (Gaussian Heuristic). *The length $\lambda_1(L)$ of the shortest vector in an n -dimensional random lattice L satisfy*

$$\lambda_1(L) \approx \sqrt{\frac{n}{2\pi e}} (\det(L))^{1/n}. \quad (3.4)$$

In particular, assume that there is a positive number C_n , depending on n , such that

$$\lambda_1(L) \geq C_n \sqrt{\frac{n}{2\pi e}} (\det(L))^{1/n}. \quad (3.5)$$

Remark 3.2.2. *Even though it seems to be a difficult problem to precisely estimate C_n in Heuristic 3.2.1.1 for all random lattices, Yuanmi and Nguyen state in Section 4.3 in [92], that the ratio between the final norm and the Gaussian heuristic prediction is mostly within 0.95 and 1.05. Therefore, one may replace C_n in Heuristic 3.2.1.1 by 0.95. It is worth noting other research papers providing estimates for C_n [68, 46], which can be used to replace C_n in Heuristic 3.2.1.1.*

Theorem 3.2.1. *Suppose that a prime order group \mathbb{G} has a t -factorization property with respect to*

$$\mathbb{B} = \{u_i \in \mathbb{G} : i = 1, \dots, c\}, \quad (3.6)$$

where $|\mathbb{B}| = c$ and u_i are chosen uniformly and independently in \mathbb{G} . Moreover, suppose that

$$|\mathbb{G}| > \frac{1}{C_{N+1}^{N+1}} (2(\sqrt{2\pi e})(b-1))^{2t+1}, \quad (3.7)$$

for some integer $b \geq 2$, for all $1 \leq N \leq 2t$, and that Heuristic 3.2.1.1 holds. Then for all $g \in \mathbb{G}_{\mathbb{B}, t}$, there is a unique ordered integer sequence $[\delta_i]_{i=1}^{|\mathbb{B}|}$ with $\sum_{i=1}^{|\mathbb{B}|} |\delta_i| \leq t$ and $|\delta_i| \leq (b-1)$ such that $g = \prod_{i=1}^{|\mathbb{B}|} u_i^{\delta_i}$.

Proof. Suppose for contradiction that there exist two distinct ordered integer sequences $[\delta_i]_{i=1}^{|\mathbb{B}|}$ and $[\tau_i]_{i=1}^{|\mathbb{B}|}$, with $\sum_{i=1}^{|\mathbb{B}|} |\delta_i| \leq t$, $|\delta_i| \leq (b-1)$, $\sum_{i=1}^{|\mathbb{B}|} |\tau_i| \leq t$, $|\tau_i| \leq (b-1)$, such that

$$\prod_{i=1}^{|\mathbb{B}|} u_i^{\delta_i} = \prod_{i=1}^{|\mathbb{B}|} u_i^{\tau_i}. \quad (3.8)$$

Moreover, suppose that $I = \{i_1, \dots, i_N\} \subseteq \{1, \dots, |\mathbb{B}|\}$ is the set of indices for which $\gamma_j = \delta_{i_j} - \tau_{i_j} \neq 0$, and that $u_{i_j} = g^{r_j}$ for some integer $r_j \in [1, |\mathbb{B}|]$, where g is a generator of $\mathbb{G} = \langle g \rangle$. Note that $1 \leq N \leq 2t$. The equation (3.8) is equivalent to

$$\sum_{j=1}^N r_j \gamma_j - k|\mathbb{G}| = 0, \quad (3.9)$$

for some integer $k \in [1, |\mathbb{G}|]$. In other words, the vector

$$\gamma = [\gamma_1, \dots, \gamma_N, 0] \quad (3.10)$$

belongs to the integer lattice L generated by the rows of the $(N+1) \times (N+1)$ matrix

$$M = \begin{bmatrix} 1 & 0 & \dots & \dots & \dots & 0 & r_1 \\ 0 & 1 & 0 & \dots & \dots & \vdots & r_2 \\ \vdots & 0 & 1 & 0 & \dots & \vdots & \vdots \\ \vdots & \vdots & 0 & \ddots & 0 & \vdots & \vdots \\ \vdots & \vdots & \vdots & 0 & 1 & 0 & r_{N-1} \\ \vdots & \vdots & \vdots & \vdots & 0 & 1 & r_N \\ 0 & \dots & \dots & \dots & \dots & 0 & |\mathbb{G}| \end{bmatrix}$$

because $\gamma = [\gamma_1, \dots, \gamma_N, -k] \times M$. The length of $\gamma \in L$ can be estimated as

$$\|\gamma\| = \sqrt{\sum_{j=1}^N \gamma_j^2} \leq 2\sqrt{N}(b-1), \quad (3.11)$$

because $|\gamma_j| = |\delta_{ij} - \tau_{ij}| \leq 2(b-1)$ for all $j = 1, \dots, N$. By assumption, u_i are uniformly and independently drawn from \mathbb{G} . Then, r_i are uniformly and independently drawn from $\{1, \dots, |\mathbb{G}| - 1\}$ where $|\mathbb{G}|$ is prime. Hence, the lattice L can be assumed random [117] and the Gaussian heuristic 3.2.1.1 implies that the length of the shortest vector in L is

$$\lambda_1(L) \geq C_{N+1} \sqrt{\frac{N+1}{2\pi e}} |\mathbb{G}|^{1/(N+1)}. \quad (3.12)$$

Finally, using the inequalities (3.11), (3.7), the approximation (3.12), and our previous observation $N \leq 2t$, we derive

$$\|v\| \leq 2\sqrt{N}(b-1) \leq 2\sqrt{N+1}(b-1) \quad (3.13)$$

$$= \sqrt{\frac{N+1}{2\pi e}} \left((2\sqrt{2\pi e}(b-1))^{N+1} \right)^{1/(N+1)} \quad (3.14)$$

$$\leq \sqrt{\frac{N+1}{2\pi e}} \left((2\sqrt{2\pi e}(b-1))^{2t+1} \right)^{1/(N+1)} \quad (3.15)$$

$$< \sqrt{\frac{N+1}{2\pi e}} (|\mathbb{G}| \cdot C_{N+1}^{N+1})^{1/(N+1)} \leq \lambda_1(L). \quad (3.16)$$

This is a contradiction because the norm of a non-zero lattice vector in L cannot be smaller than $\lambda_1(L)$. ■

Corollary 3.2.1. *Suppose that a prime order group \mathbb{G} has a t -factorization property with respect to*

$$\mathbb{B} = \{u_i \in \mathbb{G} : i = 1, \dots, c\}, \quad (3.17)$$

where $|\mathbb{B}| = c$ and u_i are chosen uniformly and independently in \mathbb{G} . Moreover, suppose that

$$|\mathbb{G}| > \frac{1}{C_{N+1}^{N+1}} (2(\sqrt{2\pi e})t)^{2t+1}, \quad (3.18)$$

for all $1 \leq N \leq 2t$, and that Heuristic 3.2.1.1 holds. Then \mathbb{G} has a unique t -factorization property with respect to \mathbb{B} .

Proof. The proof follows by replacing b in Theorem 3.2.1 by $(t+1)$. ■

Remark 3.2.3. If one uses the Chen and Nguyen's estimate $C_n = 0.95$ (see Remark 3.2.2), then the inequality (3.18) in Corollary 3.2.1 can be rewritten as

$$|\mathbb{G}| > \frac{1}{0.95^{N+1}} (2(\sqrt{2\pi e})t)^{2t+1}.$$

Similarly, if one uses the Chen's estimate $C_n = 0.147^{1/n}$, then the inequality (3.18) in Corollary 3.2.1 can be rewritten as

$$|\mathbb{G}| > \frac{1}{0.147} (2(\sqrt{2\pi e})t)^{2t+1}.$$

3.2.2 Computational Secure Sketch Construction

The main result of this section is Theorem 3.2.3, which shows that sketch schemes with respect to the L_1 distance can be constructed using groups with unique factorization. Given $x, y \in \mathbb{Z}^n$, we recall that the L_1 distance between x and y is given by

$$L_1(x, y) = \sum_{i=1}^n |x_i - y_i|.$$

We start by defining a *sketch scheme*. Even though our definition is closely related to the previous definitions of (secure) sketches, there are two significant differences to point out. First, our Definition 3.2.2 avoids associating entropy-based security notions to sketch schemes and allows us to be more flexible in constructing sketch schemes and to discuss their security based on decisional and computational problems, rather independent of the entropy of the input space conditioned on the sketch values. Second, our sketch function outputs a pair of values, where the first value explicitly enforces randomization, whereas in a traditional sketch scheme, the sketch function outputs a single (sketch) value and the randomization is built into the process.

Definition 3.2.2 (Sketch Scheme). *Let λ be a security parameter. Let t be a positive real number and \mathbb{M} a metric space with a distance function $d : \mathbb{M} \times \mathbb{M} \rightarrow \mathbb{R}$. A sketch scheme with threshold t is a tuple of randomized procedures $SS = (\text{ParamGen}, \text{Sketch}, \text{Rec})$ such that*

$$\mathbb{M}, \mathcal{R} \leftarrow \text{ParamGen}(\lambda),$$

$$\text{Sketch} : \mathbb{M} \rightarrow \mathcal{R} \times S$$

$$x \mapsto (R, s), \text{ where } R \leftarrow_s \mathcal{R},$$

$$\text{Rec} : \mathcal{R} \times S \times \mathbb{M} \rightarrow \mathbb{M}$$

$$(R, s, y) \mapsto x$$

and that, for all $x, y \in \mathbb{M}$ with $d(x, y) \leq t$, $\text{Rec}(\text{Sketch}(x), y) = x$, except with probability negligible in λ . Here, \mathcal{R} and S specify the output space of the randomized process Sketch .¹

Definition 3.2.3. We define $\mathcal{G}_{n,t} = \{(\mathbb{G}, \mathbb{B}, \mathcal{B}_n)\}$ as a family of triples, where \mathbb{G} is a multiplicative group with unique t -factorization property with respect to $\mathbb{B} = [u_i]_{i=1}^{|\mathbb{B}|}$, and $\mathcal{B}_n = [g_i]_{i=1}^n$ is an ordered sequence of pairwise distinct elements $g_i \in \mathbb{B}$.

Remark 3.2.4. Note that for a fixed \mathbb{G} , there can be multiple choices for \mathbb{B} in $\mathcal{G}_{n,t} = \{(\mathbb{G}, \mathbb{B}, \mathcal{B}_n)\}$; and for fixed \mathbb{G} and \mathbb{B} , there can be multiple choices for \mathcal{B}_n .

Definition 3.2.4 (Sketch). Let n and t be positive integers and $A \subseteq \mathbb{Z}$ a finite subset of \mathbb{Z} . We define a sketch algorithm as follows:

$$\begin{aligned} \text{Sketch}_{n,t} : A^n &\rightarrow \mathcal{G}_{n,t} \times \mathbb{G} \\ x = (x_1, \dots, x_n) &\mapsto \left(R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [g_1, \dots, g_n]), s = \prod_{i=1}^n g_i^{x_i} \right), \end{aligned}$$

where $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n) \leftarrow_s \mathcal{G}_{n,t}$.

Next, Theorem 3.2.2 shows that $\text{Sketch}_{n,t}$ can be associated with a recovery function Rec .

Theorem 3.2.2. Let $x \in A^n$ and $\text{Sketch}_{n,t}(x) = (R, s)$. Let $y \in A^n$ such that $L_1(x, y) \leq t$. Then, there exists an algorithm Rec that takes as input R, s , and y ; calls Factor as a subroutine and outputs (i.e., recovers) x .

Proof. Observe that

$$\Delta = \frac{s}{\prod_{i=1}^n g_i^{y_i}} = \prod_{i=1}^n g_i^{x_i - y_i} \in \mathbb{G}_{\mathbb{B},t}$$

because $g_i \in \mathbb{B}$ and $x_i - y_i \in \mathbb{Z}$ with $\sum_{i=1}^n |x_i - y_i| \leq t$. Therefore, Rec can compute Δ and, using a subroutine call to Factor with input Δ , it can obtain an ordered sequence of integers $[\delta_k]_{k=1}^{|\mathbb{B}|}$ such that

$$\Delta = \prod_{k=1}^{|\mathbb{B}|} u_k^{\delta_k},$$

where $\mathbb{B} = \{u_k \in \mathbb{G} : k = 1, \dots, |\mathbb{B}|\}$ and $\sum_{k=1}^{|\mathbb{B}|} |\delta_k| \leq t$. Notice that, for each g_i , there exists a unique $k_i \in \{1, \dots, |\mathbb{B}|\}$ such that $g_i = u_{k_i}$ and $x_i - y_i = \delta_{k_i}$ because \mathbb{G} has unique t -factorization property. Moreover, the (u_{k_i}, δ_{k_i}) pair can be efficiently identified from the factorization of Δ and the knowledge of \mathcal{B}_n , and hence Rec can recover x by setting $x_i = y_i + \delta_{k_i}$ for $i = 1, \dots, n$. ■

Finally, we define our subset product sketch scheme in Definition 3.2.5 and show in Theorem 3.2.3 that SPS is a sketch scheme.

Definition 3.2.5 (Subset Product Sketch (SPS)). A subset product sketch (SPS) is a triple of randomized procedures $(\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$, where ParamGen and $\text{Sketch}_{n,t}$ are defined as in Definition 3.2.4 with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$; and Rec is defined as in Theorem 3.2.2.

1. Intuitively, \mathcal{R} specifies the family of domain parameters while S specifies the domain for sketch values.

Theorem 3.2.3. $SPS = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ satisfies Definition 3.2.2 with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . That is, SPS is a sketch scheme with threshold t with respect to the L_1 distance.

Proof. The proof follows from Definition 3.2.2, Definition 3.2.4, and Theorem 3.2.2. \blacksquare

Remark 3.2.5 (Homomorphic Properties). *It is noteworthy that the sketch function exhibits homomorphic properties. It is possible to compute the product of a template and a scalar in the encrypted domain, thereby enabling the performance of blinding and unblinding operations. The result of this computation is given by the following equation:*

$$\text{Sketch}_{n,t}(kx) = \prod_{i=1}^n \left(g_i^{x_i} \right)^k = s^k \in \mathbb{G}.$$

Moreover, the addition of two vectors can be performed in the exponent to modify the value of certain coordinates if the basis is reused. The result is a product of the form:

$$\text{Sketch}_{n,t}(x + y) = \prod_{i=1}^n g_i^{x_i + y_i} = \prod_{i=1}^n g_i^{x_i} \times \prod_{i=1}^n g_i^{y_i} = \text{Sketch}_{n,t}(x) \times \text{Sketch}_{n,t}(y).$$

The aforementioned two properties show significant utility and are employed in the generation of secure protocols in this chapter (see Section 3.3).

3.2.3 Security of SPS

In this section, we show that SPS (see Definition 3.2.5) satisfies the reusability and irreversibility security properties under certain plausible assumptions.

Security Definitions and Games

Let $SS = (\text{ParamGen}, \text{Sketch}, \text{Rec})$ be a sketch scheme with threshold t as in Definition 3.2.2, \mathcal{A} a probabilistic polynomial time algorithm with access to SS and \mathcal{D} a distribution over \mathbb{M} . We say that a problem parameterized by λ is hard if the success probability of all probabilistic polynomial-time algorithms to solve that problem is negligible in λ . We first adapt reusability [61] and irreversibility [108] notions for our purposes. The concept of reusability was initially defined for fuzzy extractors [122] but, to the best of our knowledge, we are the first to adapt this property to secure sketches.

The *reusability* property describes a scenario in which an adversary has access to multiple sketches of an individual, each subject to adaptively chosen perturbations. The adversary is challenged to determine whether a new sketch belongs to the same individual.

Definition 3.2.6 (The reusability experiment $\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{REU}(\lambda)$). *Let K be a positive integer polynomial in λ . The reusability experiment $\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{REU}(\lambda)$ is defined as follows.*

1. $\mathbb{M}, \mathcal{R} \leftarrow \text{ParamGen}(\lambda)$
2. $b \leftarrow_{\$} \{0, 1\}$; $x \leftarrow_{\$} \mathcal{D}$; $(R, s) \leftarrow \text{Sketch}(x)$
3. $\psi_0 \leftarrow \perp$; $s_0 \leftarrow \perp$
4. \mathcal{A} makes K adaptive queries for $i = 1, \dots, K$

- (a) $\mathbb{M} \ni \psi_i \leftarrow \mathcal{A}(\mathbb{M}, \mathcal{R}; R, s; R_0, \dots, R_{i-1}; s_0, \dots, s_{i-1}; \psi_0, \dots, \psi_{i-1})$
- (b) $(R_i, s_i) \leftarrow \text{Sketch}(\psi_i + x)$
5. If $b = 1$, then $y \leftarrow_s \mathcal{D}$ with $d(x, y) \leq t$
6. If $b = 0$, then $y \leftarrow_s \mathcal{D}$ with $d(x, y) > t$
7. $(R', s') \leftarrow \text{Sketch}(y)$
8. $b' \leftarrow \mathcal{A}(\mathbb{M}, \mathcal{R}; R, s; R_0, \dots, R_K; s_0, \dots, s_K; \psi_0, \dots, \psi_K; R', s')$
9. If $b = b'$, then return 1; otherwise return 0

Definition 3.2.7 (Reusable Sketch). A sketch scheme $SS = (\text{ParamGen}, \text{Sketch}, \text{Rec})$ is said to be reusable with respect to the distribution \mathcal{D} if

$$\text{Adv}_{\mathcal{A}, \text{REU}}(\lambda) = \left| \mathbb{P}(\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{\text{REU}}(\lambda) = 1) - \frac{1}{2} \right|$$

is negligible in λ for all \mathcal{A} .

The *irreversibility* property models the scenario where an attacker is challenged to recover the secret input from which the sketch is derived. An adversary can simply guess by sampling a vector at random and under this naive strategy, the probability of success can be measured relative to the size of the closed ball of radius t in the n -dimensional space A^n , centered at the input vector. A successful adversary should capture better strategies, hence motivating the following definition.

Definition 3.2.8 (The Irreversibility Experiment $\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{\text{IRR}}(\lambda)$). The irreversibility experiment $\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{\text{IRR}}(\lambda)$ is defined as follows.

1. $\mathbb{M}, \mathcal{R} \leftarrow \text{ParamGen}(\lambda)$
2. $x \leftarrow \mathcal{D}; (R, s) \leftarrow \text{Sketch}(x)$
3. $y \leftarrow \mathcal{A}(\mathbb{M}, \mathcal{R}; R, s)$
4. If $d(x, y) \leq t$, then return 1; otherwise, return 0

Definition 3.2.9 (Irreversible Sketch). A sketch scheme $SS = (\text{ParamGen}, \text{Sketch}, \text{Rec})$ is said to be irreversible with respect to the distribution \mathcal{D} if

$$\text{Adv}_{\mathcal{A}, \text{IRR}}(\lambda) = \left| \mathbb{P}(\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{\text{IRR}}(\lambda) = 1) - V_t \right|,$$

is negligible for all \mathcal{A} . Here,

$$V_t = \max_{x \in \mathbb{M}} \frac{|\{y \in \mathbb{M} : d(x, y) \leq t\}|}{|\mathbb{M}|}$$

estimates the success probability of the naive \mathcal{A} returning $y \leftarrow_s \mathbb{M}$ at step (3) in $\text{Exp}_{SS, \mathcal{D}, \mathcal{A}}^{\text{IRR}}$.

Theorem 3.2.4 (Reusable implies irreversible). Let $SS = (\text{ParamGen}, \text{Sketch}, \text{Rec})$ be a sketch scheme with polynomial time Sketch and Rec algorithm. If SS is reusable with respect to \mathcal{D} , then SS is irreversible with respect to \mathcal{D} .

Proof. Suppose that SS is not irreversible. Then there exists an adversary \mathcal{A} such that $\text{Adv}_{\mathcal{A}, \text{IRR}}(\lambda)$ is non-negligible. In the following, we construct an adversary \mathcal{A}' such that $\text{Adv}_{\mathcal{A}', \text{REU}}(\lambda)$ is non-negligible. In other words, we show that SS is not reusable. \mathcal{A}' uses \mathcal{A} as a subroutine with the following strategy:

1. \mathcal{A}' skips the adaptive queries and receives $(\mathbb{M}, \mathcal{R}; R, s; R', s')$ as in $\text{Exp}_{\text{SS}, \mathcal{D}, \mathcal{A}'}^{\text{REU}}(\lambda)$
2. \mathcal{A}' runs \mathcal{A} , and gets $y \leftarrow \mathcal{A}(\mathbb{M}, \mathcal{R}; R, s)$, $y' \leftarrow \mathcal{A}(\mathbb{M}, \mathcal{R}; R', s')$
3. \mathcal{A}' runs Rec , and obtains $x \leftarrow \text{Rec}(R, s, y)$, $y \leftarrow \text{Rec}(R, s', y')$
4. \mathcal{A}' outputs $b' = 1$ if $d(x, y) \leq t$; and outputs $b' = 0$, otherwise
5. \mathcal{A}' outputs $b' \leftarrow_{\$} \{0, 1\}$ if \mathcal{A} or Rec fails

We conclude that the advantage of the adversary for reusability $\text{Adv}_{\mathcal{A}', \text{REU}}(\lambda)$ is non-negligible because $\text{Adv}_{\mathcal{A}, \text{IRR}}(\lambda)$ is non-negligible and that, for all $x, y \in \mathbb{M}$ with $d(x, y) \leq t$, $\text{Rec}(\text{Sketch}(x), y) = x$ by definition, except with negligible probability. \blacksquare

SPS is Reusable and Irreversible

In the following, we show that SPS is reusable and irreversible under the hardness assumption of the *decisional subset product problem* (DSPP) and *computational subset product problem* (CSPP), respectively. We should note that DSPP and CSPP generalize the decisional distributional modular subset product and distributional modular subset product problems as defined in [49], which consider only binary vectors and Hamming distance.

Problem 1 (Decisional Subset Product Problem (DSPP)). *Let $A \subseteq \mathbb{Z}$ and $\mathcal{B}_n = [g_i]_{i=1}^n$, $g_i \in \mathbb{G}$. Let \mathcal{D} be a distribution over A^n . Define the distribution $\mathcal{D}_0 = (\mathcal{B}_n, X)$ where $X = \prod_{i=1}^n g_i^{x_i} \in \mathbb{G}$ with $x = (x_1, \dots, x_n) \leftarrow_{\$} \mathcal{D}$. Define the distribution $\mathcal{D}_1 = (\mathcal{B}_n, X')$ where $X' \leftarrow_{\$} \mathbb{G}$. The decisional subset product problem (DSPP) in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} is to distinguish \mathcal{D}_0 from \mathcal{D}_1 . When \mathcal{D} is the uniform distribution over A^n , we call this problem \mathcal{U} -DSPP in \mathbb{G} with respect to \mathcal{B}_n .*

Theorem 3.2.5 (DSPP implies reusable). *Let $\text{SPS} = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ be a sketch scheme with $\mathbb{M} = A^n$ and $\mathcal{R} = \mathcal{G}_{n,t}$. If DSPP in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} is hard for all $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{R}$, then SPS is reusable with respect to \mathcal{D} .*

Proof. Consider the reusability experiment $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}}(\lambda)$ and let S_0 denote the event that $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}}(\lambda)$ outputs 1. Moreover, assume that in step (4b), we have $(R_i, s_i) \leftarrow \text{Sketch}_{n,t}(\psi_i + s_i)$, where $s_i = \prod_{j=1}^n g_j^{x_j + \psi_{i,j}} = \left(\prod_{j=1}^n g_j^{x_j} \right) \left(\prod_{j=1}^n g_j^{\psi_{i,j}} \right)$, which is indistinguishable from a random element in \mathbb{G} , because $\prod_{j=1}^n g_j^{x_j}$ is indistinguishable from a random element in \mathbb{G} if the DSPP is hard. Therefore, an hybrid $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}-1}(\lambda)$ can be defined by replacing s_1 by a random element of \mathbb{G} in step (4b) in $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}}(\lambda)$. Similarly, $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}-i}(\lambda)$ can be defined by replacing s_i by a random element of \mathbb{G} in step (4b) in $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}-(i-1)}(\lambda)$ for $i = 2, \dots, K$. Observe that the probability of the event S_K that $\text{Exp}_{\text{SPS}, \mathcal{A}}^{\text{REU}-K}(\lambda)$ outputs 1 is $1/2$. Using a sequence of hybrid arguments and the triangle inequality, we obtain $|\mathbb{P}(S_0) - 1/2| \leq K \times \text{Adv}_{\mathcal{A}, \text{DSPP}}(\lambda)$. \blacksquare

Corollary 3.2.2 (DSPP implies irreversible). *Let $\text{SPS} = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ be a polynomial time sketch scheme with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . If DSPP in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} is hard for all $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{G}_{n,t}$, then SPS is irreversible with respect to \mathcal{D} .*

Proof. The proof follows from Theorem 3.2.5 and Theorem 3.2.4. \blacksquare

Corollary 3.2.2 assures the irreversibility of SPS if the DSPP is hard. Next, we provide an alternative argument for the irreversibility of SPS under the hardness assumption of the discrete logarithm problem (DLP). We first recall the definition of the DLP and define the *computational subset product problem* (CSPP).

Problem 2 (Discrete logarithm problem). *The discrete logarithm problem (DLP) in \mathbb{G} with respect to g is the following: Given g and $h = g^x \in \mathbb{G}$ for some (unknown) $x \leftarrow_{\$} \mathbb{Z}_{|\mathbb{G}|}^*$, compute x .*

Problem 3 (Computational Subset Product Problem (CSPP)). *Let $A \subseteq \mathbb{Z}$ and $\mathcal{B}_n = [g_i]_{i=1}^n$, $g_i \in \mathbb{G}$. Let \mathcal{D} be a distribution over A^n . The computational subset product problem (CSPP) in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} is the following: Given $\mathcal{B}_n = [g_i]_{i=1}^n$ and $s = \prod_{i=1}^n g_i^{x_i}$ for $x = (x_1, \dots, x_n) \leftarrow_{\$} \mathcal{D}$, compute $y = (y_1, \dots, y_n) \in \mathbb{Z}^n$ such that $s = \prod_{i=1}^n g_i^{y_i}$. When \mathcal{D} is the uniform distribution over A^n , we call this problem \mathcal{U} -CSPP in \mathbb{G} with respect to \mathcal{B}_n .*

Theorem 3.2.6 (CSPP implies irreversible). *Let $\text{SPS} = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ be a polynomial time sketch scheme with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . If CSPP in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} is hard for all $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{G}_{n,t}$, then SPS is irreversible with respect to \mathcal{D} .*

Proof. Let $s = \prod_{i=1}^n g_i^{x_i}$ be a given instance of the CSPP with $\mathcal{B}_n = [g_i]_{i=1}^n$ and $x = (x_1, \dots, x_n) \leftarrow_{\$} \mathcal{D}$ for some distribution \mathcal{D} over A^n . Suppose that SPS is not irreversible. Then there exists an adversary \mathcal{A} such that $\text{Adv}_{\mathcal{A}, \text{IRR}}(\lambda)$ is non-negligible. In other words, \mathcal{A} can output $y \in A^n$ such that $d(x, y) \leq t$. Now, given y and s , the recovery algorithm Rec can output x in polynomial time with a non-negligible probability. Hence, CSPP can be solved in polynomial time with a non-negligible probability and that finishes the proof. \blacksquare

3.2.4 On the Hardness of and Relationships between DSPP and CSPP

As discussed previously, reusability and irreversibility of SPS rely on the hardness of DSPP and CSPP, respectively. In this section, we study the hardness of DSPP, CSPP, and study some relationships between these problems and the DLP. We first show that the hardness of DLP and \mathcal{U} -DSPP implies the hardness of \mathcal{U} -CSPP (Theorem 3.2.7). This can be seen as evidence that CSPP may be harder than DSPP. We also show that the hardness of DLP implies the hardness of CSPP if the underlying Sketch in SPS is surjective (Theorem 3.2.8). Here, our motivation to introduce surjectivity for Sketch is to relate the subset product problems and the security of SPS to other well-known problems in cryptography. Based on Theorem 3.2.7 and Theorem 3.2.8, it is natural to ask if there is a strong relationship between the hardness of DSPP and surjectivity of Sketch. Figure 3.1 provides a summary of our results.

Theorem 3.2.7 (\mathcal{U} -DSPP and DLP implies \mathcal{U} -CSPP). *Let $\text{SPS} = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ be a sketch scheme with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . Suppose that $\mathbb{G} = \langle g \rangle$ is a cyclic group generated by g and that \mathcal{U} -DSPP with respect to \mathcal{B}_n is hard. If there is an algorithm that solves \mathcal{U} -CSPP in \mathbb{G} with respect to \mathcal{B}_n in time $T_{\mathcal{U}\text{-CSPP}}$ for all $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{G}_{n,t}$, then DLP in \mathbb{G} with respect to g can be solved in time $\mathcal{O}(nT_{\mathcal{U}\text{-CSPP}})$.*

Proof. Let $h \in \mathbb{G}$ be given as an instance of the DLP, and let $\mathcal{A}_{\mathcal{U}\text{-CSPP}}$ be an algorithm that solves \mathcal{U} -CSPP in time $T_{\mathcal{U}\text{-CSPP}}$. We describe an algorithm \mathcal{A} that computes $a \in \mathbb{Z}$ such that $h = g^a$. First, \mathcal{A} computes $s_k = g^{a_k}$ for $a_k \leftarrow_{\$} \mathbb{Z}_{|\mathbb{G}|}$, and calls $\mathcal{A}_{\mathcal{U}\text{-CSPP}}$ with input s_k and $\mathcal{B}_n = [g_i]_{i=1}^n$. Since \mathcal{U} -DSPP is hard,

$$s_k = \prod_{i=1}^n g_i^{x_{k,i}}, \quad x_{k,i} \in \mathbb{Z}. \quad (3.19)$$

for some $(x_{k,1}, \dots, x_{k,n}) \leftarrow_s A^n$, with non-negligible probability. Hence, \mathcal{A} will receive $\mathbf{x}_k = (x_{k,1}, \dots, x_{k,n})$ as output of $\mathcal{A}_{\mathcal{U}\text{-CSPP}}$. As a result, \mathcal{A} obtains a modular linear relation

$$a_k \equiv \sum_{i=1}^n x_{k,i} d_i \pmod{|\mathbb{G}|}, \quad (3.20)$$

where $g_i = g^{d_i}$ for some integers d_i and $i = 1, \dots, n$. \mathcal{A} repeats this process until it obtains n linearly independent relations, where the total number of repetitions is expected to be polynomial in n . After collecting n linearly independent relations, \mathcal{A} can recover d_i for all $i = 1, \dots, n$ by solving a linear system of equations in time polynomial in n . Finally, \mathcal{A} computes $s = h^b = g^{ab}$ for some random integer b relatively prime with $|\mathbb{G}|$; calls $\mathcal{A}_{\mathcal{U}\text{-CSPP}}$ with input s and \mathcal{B}_n , until receiving $\mathbf{x} = (x_1, \dots, x_n)$, $x_i \in \mathbb{Z}$, such that

$$s = \prod_{i=1}^n g_i^{x_i}, \quad (3.21)$$

and recovers the discrete logarithm a of h with respect to g via the modular equation

$$a \equiv b^{-1} \left(\sum_{i=1}^n x_i d_i \right) \pmod{|\mathbb{G}|}. \quad (3.22)$$

■

Definition 3.2.10 (Surjective Sketch). *We say that a Sketch $: \mathbb{M} \rightarrow \mathcal{R} \times S$ is surjective if for any $R \in \mathcal{R}$ and $s \in S$, there exists $\mathbf{x} \in \mathbb{M}$ such that Sketch(\mathbf{x}) = (R, s).*

Example 3.2.4.1. Let $\mathbb{G} = \mathbb{Z}_{13}^*$, $\mathbb{B} = \{2, 3, 5\}$, and $t = 1$. Notice that \mathbb{G} has a unique t -factorization property with respect to \mathbb{B} . Let $n = 2$, $A = \{0, 1, 2, 3\}$, and

$$\mathcal{G}_{n,t} = \{(\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [2, 3]), (\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [2, 5]), (\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [3, 5])\}.$$

One can easily verify that Sketch $: A^n \rightarrow \mathcal{G}_{n,t} \times \mathbb{G}$ is surjective.

Theorem 3.2.8. (Surjective Sketch and DLP implies CSPP) *Let SPS = (ParamGen, Sketch $_{n,t}$, Rec) be a sketch scheme with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . Suppose that Sketch $_{n,t}$ is surjective, and $\mathbb{G} = \langle g \rangle$ is a cyclic group generated by g . If there is an algorithm that solves*

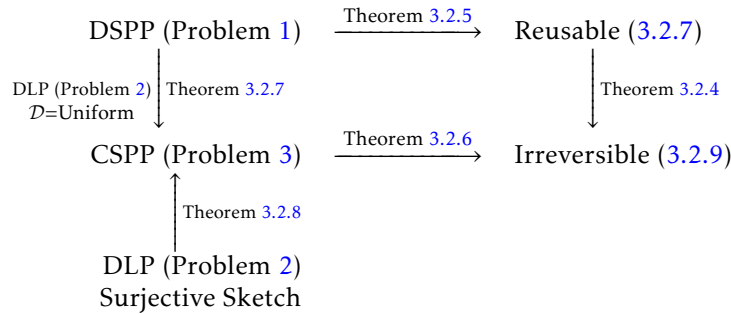


Figure 3.1 – Relations between the hardness of problems DSPP, CSPP, DLP, and the security properties reusability and irreversibility of SPS.

CSPP in \mathbb{G} with respect to \mathcal{B}_n and \mathcal{D} in time T_{CSPP} for all $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{G}_{n,t}$, then DLP in \mathbb{G} with respect to g can be solved in time $\mathcal{O}(nT_{CSPP})$.

Proof. Let $h \in \mathbb{G}$ be given as an instance of the DLP, and let \mathcal{A}_{CSPP} be an algorithm that solves CSPP in time T_{CSPP} . We describe an algorithm \mathcal{A} that computes $a \in \mathbb{Z}$ such that $h = g^a$. First, \mathcal{A} computes $s_k = g^{a_k}$ for randomly chosen integers $a_k \in \mathbb{Z}_{|\mathbb{G}|}$ and calls \mathcal{A}_{CSPP} with input s_k and $\mathcal{B}_n = [g_i]_{i=1}^n$. Since $\text{Sketch}_{n,t}$ is surjective, \mathcal{A} will receive $x_k = (x_{k,1}, \dots, x_{k,n})$ as output of \mathcal{A}_{CSPP} , where

$$s_k = \prod_{i=1}^n g_i^{x_{k,i}}, \quad x_{k,i} \in \mathbb{Z}. \quad (3.23)$$

The rest of the proof follows similarly the proof of Theorem 3.2.7. ■

Corollary 3.2.3. (*Surjective Sketch and DLP implies irreversible*)

Let $SPS = (\text{ParamGen}, \text{Sketch}_{n,t}, \text{Rec})$ be a sketch scheme with $\mathbb{M} = A^n$, $\mathcal{R} = \mathcal{G}_{n,t}$, $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n)$, $S = \mathbb{G}$, $d = L_1$, and a threshold t . Suppose that $\mathbb{G} = \langle g \rangle$ is cyclic and $\text{Sketch}_{n,t}$ is surjective for all $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \in \mathcal{G}_{n,t}$. If DLP in \mathbb{G} with respect to g is hard, then SPS is irreversible.

Proof. The proof follows from Theorem 3.2.6 and Theorem 3.2.8. ■

Estimating the Size of Balls

We denote the set of integers and real numbers by \mathbb{Z} and \mathbb{R} , respectively. For positive integers $b \geq 2$ and n , we define

$$A_b = \{0, 1, \dots, b-1\},$$

as the set of integers from 0 to $b-1$; and

$$A_b^n = \{x = (x_1, \dots, x_n) : x_i \in A_b\},$$

as the set of length- n vectors over A_b .

In this chapter, we are interested in two different types of distance functions

$$d : A_b^n \times A_b^n \rightarrow \mathbb{R}$$

on A_b^n , namely the Hamming distance HD, and the Manhattan distance L_1 , which are defined as follows:

$$\begin{aligned} d(x, y) &= \text{HD}(x, y) = \#\{i : x_i \neq y_i\} \\ d(x, y) &= L_1(x, y) = \sum_{i=1}^n |x_i - y_i|, \end{aligned}$$

where $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ are in A_b^n .

Notice that, when $b = 2$, HD and L_1 are equivalent. For given $x \in A_b^n$, $t \in \mathbb{R}$, and a distance function d on A_b^n , the ball of radius t about its center x with respect to d , denoted $\text{Ball}_{A_b^n, d}(x, t)$, is defined as

$$\text{Ball}_{A_b^n, d}(x, t) = \{y \in A_b^n : d(x, y) \leq t\}.$$

When $d = \text{HD}$, the size of the ball $\text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t)$ is independent of its center \mathbf{x} , and which leads to the definition of the *volume of a Hamming ball of radius t* [20]:

$$\text{Vol}_{A_b^n, \text{HD}}(t) = \left| \text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t) \right| = \sum_{k=0}^t \binom{n}{k} (b-1)^k \quad (3.24)$$

Some of our security discussions in this chapter rely on estimating $\left| \text{Ball}_{A_b^n, d}(\mathbf{x}, t) \right|$. When $d = \text{HD}$, this boils down to estimating $\text{Vol}_{A_b^n, \text{HD}}(t)$ (see (3.24)). for which we refer to the definition of the b -ary entropy function $H_b(\alpha)$ (Definition 3.2.11) and to the well-known result Theorem 3.2.9 (a proof can be found in Section 3.3.1 in [20]).

Definition 3.2.11 (b -ary entropy function). *For an integer $b \geq 2$ and $0 \leq \alpha \leq 1$, the b -ary entropy function is defined as*

$$H_b(\alpha) = \alpha \log_b(b-1) - \alpha \log_b(\alpha) - (1-\alpha) \log_b(1-\alpha)$$

Theorem 3.2.9 (Estimating $\text{Vol}_{A_b^n, \text{HD}}(t) = \left| \text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t) \right|$). *Let $b \geq 2$ be an integer and $0 \leq \alpha \leq (b-1)/b$. Then, for all $\mathbf{x} \in A_b^n$ and sufficiently large n , we have*

$$\begin{aligned} \left| \text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t = \alpha \cdot n) \right| &\leq b^{H_b(\alpha)n}, \\ \left| \text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t = \alpha \cdot n) \right| &\geq b^{H_b(\alpha)n - o(n)}. \end{aligned}$$

Note that for a fixed $b \geq 2$ and sufficiently large n , Theorem 3.2.9 yields the estimate

$$\left| \text{Ball}_{A_b^n, \text{HD}}(\mathbf{x}, t = \alpha \cdot n) \right| \sim b^{H_b(\alpha)n},$$

where $0 \leq \alpha \leq (b-1)/b$.

We should emphasize that, for general d , $\left| \text{Ball}_{A_b^n, d}(\mathbf{x}, t) \right|$ depends on the center \mathbf{x} of the ball, and so $\text{Vol}_{A_b^n, \text{HD}}(t)$ may not be generalized for other distance functions. In particular, we are not aware of analogous estimates for $\left| \text{Ball}_{A_b^n, d}(\mathbf{x}, t) \right|$ for a general distance function d . However, as we show in Theorem 3.2.10, we can explicitly derive computable upper and lower bounds for $\left| \text{Ball}_{A_b^n, L_1}(\mathbf{x}, t) \right|$. Theorem 3.2.10 relies on Lemma 3.2.1 from [111].

Lemma 3.2.1. [Lemma 1.1 in [111]] *Let m, k, b be integers such that $1 \leq m \leq k$, $b \geq 2$. Let $\omega_{m,b}(k)$ be the number of ordered partitions of the integer k into m parts of size between 0 and $(b-1)$. Then,*

$$\omega_{m,b}(k) = \sum_{i=0, b, 2b, \dots}^k (-1)^{i/b} \binom{m}{i/b} \binom{k-i+m-1}{k-i}. \quad (3.25)$$

For a better presentation of the proof of Theorem 3.2.10, we define the following sets and prove some results.

For $x \in A_b^n$, and positive integers k and m with $m \leq k$, define

$$\begin{aligned} S_{x,b,k} &= \{y \in A_b^n : \sum_{i=1}^n |y_i - x_i| = k\}, \\ U_{b,k} &= \{z \in \mathbb{Z}^n : 1-b \leq z_i \leq b-1, \sum_{i=1}^n |z_i| = k\}, \\ U_{b,k,m} &= \{z \in U_{b,k} : |\text{Supp}(z)| = m\}, \\ U_{b,k,m}^{\geq 0} &= \{z \in U_{b,k,m} : z_i \geq 0\}, \\ V_{b,k,m} &= \{z \in \mathbb{Z}^m : 1 \leq z_i \leq b-1, \sum_{i=1}^m z_i = k\}, \\ W_{b,k,m} &= \{z \in \mathbb{Z}^m : 0 \leq z_i \leq b-2, \sum_{i=1}^m z_i = k-m\}, \\ L_{b',k} &= \{z \in A_{b'}^n : \sum_{i=1}^n z_i = k\} \end{aligned}$$

Here, the *support* of a vector z , $\text{Supp}(z)$, is defined as the set of indices i , where the components z_i of z are non-zero. That is,

$$\text{Supp}(z) = \{i : z_i \neq 0\}.$$

Lemma 3.2.2. *Let $b \geq 2$ and t be positive integers. Then*

$$\left| \text{Ball}_{A_b^n, L_1}(x, t) \right| - 1 = \sum_{k=1}^t |S_{x,b,k}| \leq \sum_{k=1}^t |U_{b,k}| = \sum_{k=1}^t \sum_{m=1}^k |U_{b,k,m}|,$$

for all $x \in A_b^n$.

Proof. The proof easily follows by the definitions of the underlying sets, and by observing that the function

$$\begin{aligned} \phi : S_{x,b,k} &\rightarrow U_{b,k} \\ (y_1, \dots, y_n) &\mapsto (z_1, \dots, z_n), \end{aligned}$$

where $z_i = y_i - x_i$, is well-defined and injective. ■

Lemma 3.2.3. *Let $b \geq 2$, k , m be positive integers with $m \leq k$. Then*

$$|U_{b,k,m}| = 2^m |U_{b,k,m}^{\geq 0}|$$

Proof. Let ϕ be the function defined as

$$\begin{aligned} \phi : U_{b,k,m} &\rightarrow U_{b,k,m}^{\geq 0} \\ (z_1, \dots, z_n) &\mapsto (|z_1|, \dots, |z_n|). \end{aligned}$$

The proof follows because ϕ is an onto function and each element in $U_{b,k,m}^{\geq 0}$ has exactly m elements in its support and hence has exactly 2^m preimages under ϕ . ■

Lemma 3.2.4. *Let $b \geq 2$, k, m, n be positive integers with $m \leq k$ and $m \leq n$. Then*

$$|U_{b,k,m}^{\geq 0}| \leq \binom{n}{m} |V_{b,k,m}|$$

Proof. Let ϕ be the function defined as

$$\begin{aligned} \phi : U_{b,k,m}^{\geq 0} &\rightarrow V_{b,k,m} \\ (z_1, \dots, z_n) &\mapsto (z_{i_1}, \dots, z_{i_m}), \end{aligned}$$

where $\{i_j : j = 1, \dots, m\} = \text{Supp}(z)$. The proof follows because ϕ is an onto function and each element in $V_{b,k,m}$ has at most $\binom{n}{m}$ preimages under ϕ . ■

Lemma 3.2.5. *Let $b \geq 2$, k, m be positive integers with $m \leq k$. Then*

$$|V_{b,k,m}| = |W_{b,k,m}| = \omega_{m,b-1}(k-m)$$

Proof. The function

$$\begin{aligned} \phi : V_{b,k,m} &\rightarrow W_{b,k,m} \\ (z_1, \dots, z_m) &\mapsto (z_1 - 1, \dots, z_m - 1), \end{aligned}$$

is a bijection and so $|V_{b,k,m}| = |W_{b,k,m}|$. It follows that $|W_{b,k,m}| = \omega_{m,b-1}(k-m)$ because, by definition, $\omega_{m,b-1}(k-m)$ is the number of ordered partitions of the integer $(k-m)$ into m parts of size between 0 and $(b-2)$, which is precisely $|W_{b,k,m}|$. ■

Lemma 3.2.6. *Let $b \geq 2$, k, n be positive integers, and $b' = \lfloor (b-1)/2 \rfloor + 1$. Then*

$$|S_{x,b,k}| \geq |L_{b',k}| = \omega_{n,b'}(k),$$

for all $x \in A_b^n$.

Proof. Consider the function

$$\begin{aligned} \phi : L_{b',k} &\rightarrow S_{x,b,k} \\ (z_1, \dots, z_n) &\mapsto (y_1, \dots, y_n), \end{aligned}$$

where $y_i = x_i - z_i$ if $x_i > b' - 1$; and $y_i = x_i + z_i$ if $x_i \leq b' - 1$. Note that this is a well-defined function because $0 \leq y_i \leq b - 1$ and

$$\sum_{i=1}^n |y_i - x_i| = \sum_{i=1}^n |z_i| = k.$$

Also note that ϕ is injective, and so $|S_{x,b,k}| \geq |L_{b',k}|$. Finally, $|L_{b',k}| = \omega_{n,b'}(k)$ follows because, by definition, $\omega_{n,b'}(k)$ is the number of ordered partitions of the integer k into n parts of size between 0 and $b' - 1$, which is precisely $|L_{b',k}|$. ■

Theorem 3.2.10 (Estimating $|\text{Ball}_{A_b^n, L_1}(x, t)|$). *Let $b \geq 2$ be an integer, $b' = \lfloor (b-1)/2 \rfloor + 1$, $x \in A_b^n$, and $\omega_{m,b}(k)$ defined as in Lemma 3.2.1. We have*

$$|\text{Ball}_{A_b^n, L_1}(x, t)| \leq 1 + \sum_{k=1}^t \sum_{m=1}^k 2^m \binom{n}{m} \omega_{m,b-1}(k-m)$$

and

$$|\text{Ball}_{A_b^n, L_1}(x, t)| \geq 1 + \sum_{k=1}^t \omega_{n,b'}(k)$$

Proof. The upper bound follows from Lemmas 3.2.2-3.2.5. The lower bound follows from Lemma 3.2.2 and Lemma 3.2.6. ■

Theorem 3.2.11. *Let $x \in A_b^n$ and $b \geq 3$. Then,*

$$|\text{Ball}_{A_b^n, L_1}(x, t)| \leq 2^{n+2t-1}$$

Proof. Using Theorem 3.2.10, Vandermonde's identity and the bound on $\omega_{m,b-1}(k-m)$ given by Ott *et al.* [37], we have

$$\begin{aligned} |\text{Ball}_{A_b^n, L_1}(x, t)| &\leq 1 + \sum_{k=1}^t \sum_{m=1}^k 2^m \binom{n}{m} \omega_{m,b-1}(k-m) \\ &\leq \sum_{k=1}^t 2^k \sum_{m=1}^k \binom{n}{m} \binom{k-1}{m-1} \\ &\leq 2^t \sum_{k=1}^t \binom{n+k-1}{k} \leq 2^t \sum_{k=1}^t \binom{n+t-1}{k} \\ &\leq 2^{n+2t-1} \end{aligned}$$

and the result follows. ■

Security Against Known Attacks

As discussed in Section 3.2.3 and summarized in Figure 3.1, DSPP and CSPP are the main hardness problems to claim reusability and irreversibility of SPS. In addition, reusability implies irreversibility, and that irreversibility follows mainly from the hardness of DLP. As much as these reductionist arguments provide some security assurance for SPS, one should carefully study all the underlying assumptions and try to estimate the concrete security of SPS. More precisely, DSPP and CSPP assumptions depend on the choice of the basis \mathcal{B}_n and the distribution \mathcal{D} over the input space. The sketch function is assumed to be surjective when reducing DLP to the irreversibility of SPS in Theorem 3.2.3. As a worst-case scenario, the hardness of DSPP, CSPP, and DLP may fail and the sketch function may not be surjective due to the choice of \mathcal{B}_n , \mathcal{D} , and other parameters such as n , t , and \mathbb{G} . Even though some of these failures may not imply an immediate threat for the security of SPS, reductionist arguments would be inconclusive. Therefore, in this section, we follow a common practice in cryptography and try to estimate the security of SPS based on the best-known attack strategies.

We first investigate attacks on the irreversibility of SPS with threshold t with respect to the L_1 distance. Suppose that an adversary \mathcal{A} knows the parameters of SPS, namely $\mathbb{M}, \mathbb{G}, \mathbb{B}, t$, and $\mathcal{B}_n = [g_i]_{i=1}^n$. For concreteness, we furthermore assume $\mathbb{M} = A_b^n$, where $A_b = \{0, 1, \dots, b-1\}$ for some integer $b \geq 2$, and that \mathcal{D} is some distribution over A_b^n . Now, suppose that \mathcal{A} captures $X = \prod_{i=1}^n g_i^{x_i} \in \mathbb{G} = \langle g \rangle$ for some unknown $x = (x_1, \dots, x_n) \leftarrow \mathcal{D}$, and aims to output $y \in A_b^n$ such that $L_1(x, y) \leq t$. \mathcal{A} may follow the strategies as described below.

Exploit \mathcal{D} . In practice, we may not have control over the choice of \mathcal{D} . For example, x may be the encoding of a biometric input and could induce a low entropy on the input space for several reasons, such as a high correlation on the components of x . Therefore, we assume that \mathcal{A} can fully exploit \mathcal{D} and succeed in her attack with complexity

$$C_{\mathcal{D}} \approx 2^{\mu_{\mathcal{D}}}$$

for some $\mu_{\mathcal{D}} > 0$. Note that $\mu_{\mathcal{D}} = (\log_2 b)n$ would correspond to the uniform distribution \mathcal{D} over A_b^n . In practice, $\mu_{\mathcal{D}}$ is expected to be lower than $(\log_2 b)n$ and estimating μ is an active area of research.

Remark 3.2.6. An interesting way to compute $\mu_{\mathcal{D}}$ is to use the min entropy [106]. Once \mathcal{D} fully defined, all element i of the space may be associated with a probability p_i . Then, $\mu_{\mathcal{D}} = -\log_2(\max_i p_i)$. To approximate the complexity of the attacks in the non-uniform cases, a solution would be to replace the template size n by $m = \mu_{\mathcal{D}}$ in their respective complexities.

Remark 3.2.7. In the context of biometrics, Daugman's study demonstrates that 2048-bit iriscodes exhibit entropies of 249 bits [104].

Guess y . In this strategy, \mathcal{A} chooses y uniformly at random from A_b^n and hopes that $L_1(x, y) \leq t$. The success probability of this attack can be estimated as

$$\frac{|\text{Ball}_{A_b^n, L_1}(x, t)|}{|A_b^n|},$$

with $\text{Ball}_{A_b^n, L_1}(x, t) = \{y \in A_b^n : L_1(x, y) \leq t\}$ and so the complexity of the attack can be estimated as

$$C_{\text{Guess}} \approx \frac{b^n}{|\text{Ball}_{A_b^n, L_1}(x, t)|}$$

There are two cases to consider: $b = 2$ and $b \geq 3$. Estimating $|\text{Ball}_{A_b^n, L_1}(x, t)|$ in the binary case for $b = 2$, when L_1 is the same as Hamming distance, is a well-studied problem in the literature, and we have

$$|\text{Ball}_{A_b^n, L_1}(x, t)| = \sum_{k=0}^t \binom{n}{k}$$

Estimating $|\text{Ball}_{A_b^n, L_1}(x, t)|$ for $b \geq 3$ seems to be a harder problem mainly because the size of the balls is not independent of the choice of their centers. We are not aware of any previous work on this topic, and we study this problem in Section 3.2.4, which could be of independent interest. In

| n | t | $b = 2$ | | $b = 4$ | | | $b = 8$ | | | $b \in \{2, 4, 8\}$ C_S |
|------|-----|---------|--------------------|---------|---------|--------------------|---------|---------|--------------------|------------------------------|
| | | μ_G | C_{Guess} | t | μ_G | C_{Guess} | t | μ_G | C_{Guess} | |
| 128 | 16 | 0.48 | 2^{61} | 48 | 0.58 | 2^{74} | 112 | 0.66 | 2^{84} | 2^{25} |
| 256 | 32 | 0.47 | 2^{120} | 96 | 0.57 | 2^{145} | 224 | 0.64 | 2^{165} | 2^{51} |
| 640 | 80 | 0.46 | 2^{296} | 240 | 0.55 | 2^{356} | 560 | 0.63 | 2^{405} | 2^{128} |
| 1024 | 128 | 0.46 | 2^{471} | 384 | 0.55 | 2^{568} | 896 | 0.63 | 2^{645} | 2^{204} |

Table 3.1 – Concrete security estimates for the cases $n = 128, 256$, $b = 2, 4, 8$, and $t = (b-1)n/8$. In this table, C_{Guess} estimates the complexity of the guessing attack as $2^{\mu_G n}$ while C_S estimates the complexity of the attack based on solving DLP, Knapsack, and SVP as $2^{0.2n}$.

particular, in Theorem 3.2.10, we prove that

$$\left| \text{Ball}_{A_{b,L_1}^n}(x, t) \right| \leq 1 + \sum_{k=1}^t \sum_{m=1}^k 2^m \binom{n}{m} \omega_{m,b-1}(k-m),$$

where $\omega_{m,b-1}(k-m)$ is the number of ordered partitions of the integer $(k-m)$ into m parts of size in $[0, b-2]$, which can be explicitly computed using Lemma 3.2.1. As a result, we can estimate

$$C_{\text{Guess}} \gtrsim 2^{\mu_G n},$$

where $\mu_G = (\log_2 b)(1 - c_t)$ for some $c_t \geq 0$ such that

$$c_t \approx \begin{cases} \log_2 \left(\sum_{k=0}^t \binom{n}{k} \right) / n, & \text{for } b = 2, \\ \log_b \left(1 + \sum_{k=1}^t \sum_{m=1}^k 2^m \binom{n}{m} \omega_{m,b-1}(k-m) \right) / n, & \text{for } b \geq 3. \end{cases}$$

In Table 3.1, we present some concrete estimates for the cases $n = 128, 256$, $b = 2, 4, 8$, and $t = (b-1)n/8$, which corresponds to a capability of correcting a linear fraction of errors.

More generally, we prove in Theorem 3.2.12 that C_{Guess} is exponential while a linear fraction of errors can be recovered.

Theorem 3.2.12. *Let n and $b \geq 2$ be positive integers. Let $t = \alpha(b-1)n$, where $\alpha \in (0, 1/2)$ if $b = 2$, and $\alpha \in (0, (\log_2 b - 1)/(2(b-1)))$ if $b \geq 3$. Then, $C_{\text{Guess}} \gtrsim 2^{\mu_G n}$ for some $\mu_G > 0$. In other words, there exist parameters for SPS that allow recovering a linear fraction of errors (namely, up to $n/2$ errors when $b = 2$, and up to $(\log_2 b - 1)n/2$ errors when $b \geq 3$) while the complexity of the guessing attack is exponential.*

Proof. First assume $b = 2$, $t = \alpha n$, and $\alpha \in (0, 1/2)$. We can observe, using Theorem 3.2.9, that

$$C_{\text{Guess}} \gtrsim 2^{(1-H_2(\alpha))n}$$

and finish the proof by setting $\mu_G = (1 - H_2(\alpha))$ because $H_2(\alpha) < 1$ for $\alpha < 1/2$. Now, assume $b \geq 3$, $t = \alpha(b-1)n$, and $\alpha \in (0, (\log_2 b - 1)/(2(b-1)))$. We can observe, using Theorem 3.2.11, that

$$C_{\text{Guess}} \gtrsim 2^{(\log_2 b - (1 + (2t/n) - (1/n)))n}$$

and finish the proof by setting $\mu_G = (\log_2 b - (1 + (2t/n) - (1/n)))$ because one can show after some algebra that $\alpha < (\log_2 b - 1)/(2(b-1))$ implies $\mu_G > 0$. ■

Solve DLP/Knapsack/SVP. In this strategy, \mathcal{A} mounts a more sophisticated attack and first solves discrete logarithms of X and g_i for $i = 1, \dots, n$, namely $r, d_i \in \mathbb{Z}$ such that $X = g^r$ and $g_i = g^{d_i}$. This yields a modular equation

$$\sum_{i=1}^n x_i d_i \equiv r \pmod{|\mathbb{G}|}.$$

\mathcal{A} can try to solve for x_i from this equation via solving the (modular) knapsack problem (KP) or via solving the shortest vector problem (SVP) (see the proof of Theorem 3.2.1). The complexity of solving DLP is subexponential in $\log_2 |\mathbb{G}|$ and can be estimated based on the best-known attacks with respect to the characteristic of the finite field (small/medium/large characteristic) [39]. The complexity of solving SVP and KP can be estimated as $2^{0.2n}$ [41] and $2^{0.241n}$ [85], respectively. Therefore, we estimate the complexity of this attack strategy as

$$C_S \gtrsim 2^{0.2n},$$

and present some estimates for a certain set of parameters in Table 3.1.

Remark 3.2.8. (Complexity of attacking SPS is exponential in n) We are not aware of better strategies to attack the irreversibility of SPS other than the ones we discussed above. Similarly, the best approach to compromise the reusability of SPS appears to be attacking irreversibility. Hence, our analysis indicates that the complexity of attacking SPS is $2^{\mu n}$, where $\mu = \min(\mu_D/n, \mu_G, 0.2)$. Note that this complexity is exponential in n if the input space has sufficient entropy and t is chosen carefully, as explicitly described in Theorem 3.2.12.

3.2.5 Robust SPS in the Random Oracle Model

Robustness is a fundamental property in cryptography, especially when constructing a mutual authentication scheme. It states that any alteration to public information results in the primitive or protocol aborting. It is important to note that the protocol may abort due to either incorrect input or incorrect public data. Informally, robustness is a property that ensures the public data used is correct and unaltered, and that the recovery phase has been successful. The main advantage of robustness is that a user can deduce if the server that sent the public information is corrupt or malicious. This allows the user to stop all interactions before revealing any information to the server by further communicating with it.

Definition 3.2.12 (Robustness (Sketch)). Let S (Sketch, Rec) be a secure sketch for a threshold t and $Pub = S.Sketch(w)$ with $w \in \mathbb{M}$. S is said robust iff for any w' , $S.Rec(Pub', w')$ with $Pub' \neq Pub$ aborts with an overwhelming probability.

Boyer *et. al* [120] show in the random oracle model a generic way to turn any secure sketch into a robust secure sketch. Considering H a hash function that is indistinguishable from a random oracle and a non-robust secure sketch (Sketch*, Rec*), a robust secure sketch (Sketch*, Rec*) can be constructed on top of these building blocks in the following way:

| | |
|------------------------|--------------------------|
| Sketch(w) | Rec(y, pub = (pub*, h)) |
| pub* ← Sketch*(w) | w' ← Rec*(y, pub*) |
| h ← H(w, pub*) | If H(w', pub*) ≠ h abort |
| Return pub = (pub*, h) | Else, return w' |

This generic construction requires a hash function indistinguishable from a random oracle, *e.g.*, SHA3 [75]. Following the above description, our SPS can be easily modified to achieve the robustness property.

3.2.6 Concrete Instantiations of SPS

We observe that SPS can be realized in practice using the virtual black box (VBB) [49] and noise tolerant template (NTT) [71] primitives. In this section we give hints on how it is implemented and its performances. More details on the implementation are provided in Section 3.2.7.

Instantiation of $\text{Sketch}_{n,t}$. Both VBB and NTT parameter generations take as input n and t , and output a group \mathbb{G} as well as a basis $\mathcal{B}_n = [g_1, \dots, g_n]$. In the case of NTT, \mathbb{G} is a subgroup of the multiplicative group of a finite field \mathbb{F}_{q^2} , and \mathcal{B}_n consists of elements represented by the base field \mathbb{F}_p of \mathbb{F}_{q^2} . In the case of VBB, \mathbb{G} is a subgroup of the multiplicative group of integers modulo a prime q , and \mathcal{B}_n consists of small prime numbers. In both cases, \mathbb{G} and \mathcal{B}_n are used to map a binary vector $\mathbf{x} = (x_1, \dots, x_n)$ to a value $X = \prod_{i=1}^n g_i^{x_i} \in \mathbb{G}$. The transformation is referred to as *project* in NTT, and as *encode* in VBB. It is straightforward to generalize this transformation from binary vectors to \mathbf{x} with $0 \leq x_i \leq b-1$, which we use in our instantiation.

Instantiation of Rec. Both VBB and NTT propose algorithms to reconstruct the vector \mathbf{x} given X and another vector \mathbf{y} , where \mathbf{x} and \mathbf{y} are binary vectors with $\text{HD}(\mathbf{x}, \mathbf{y}) \leq t$. The reconstruction can be generalized from binary vectors to \mathbf{x}, \mathbf{y} with $0 \leq x_i, y_i \leq b-1$ when $L_1(\mathbf{x}, \mathbf{y}) \leq t$.

Reconstruction algorithms are referred to as *Decomp* in NTT, and as *Decoding* in VBB. Hence, by Theorem 3.2.1, Corollary 3.2.1, and Theorem 3.2.3, SPS can be realized using VBB and NTT under Heuristic 3.2.1.1. More details are provided in Section 3.2.7.

Performance and Security Evaluations. We provide running time evaluations for our constructions based on single-thread C and C++ programs using the GMP [13] and NTL [16] libraries.

The aforementioned processes are executed on a computer running Debian 11, which is equipped with an 11th-generation Intel Core i7-1185G7 processor operating at 3.00 GHz and 16 GB of RAM.

Table 3.2 summarizes the performance tests over 100 iterations for $n = 640$, $b \in \{2, 4, 8\}$, and $t = (b-1)n/8$ including the median as well as the standard deviation denoted by SD. Note that the parameter set provides 128-bit security level according to Table 3.1. Our implementation demonstrates that both realizations of the SPS (SPS_{NTT} and SPS_{VBB}) are efficient and suitable for applications in practice. To reduce the execution time, the parameter generation of the Sketch function can be pre-calculated.

3.2.7 Details on VBB and NTT

In this section, we provide details on our implementation of SPS using NTT [49] and VBB [71].

Instantiation of SPS using VBB

VBB [49] is an obfuscator for Hamming distance-based fuzzy matching where the security relies on the difficulty of the Modular Subset Product Problem (MSP), which can be realized as a special case of CSPP for $A = \{0, 1\}$ and $\mathbb{G} = \mathbb{Z}_q$; see Problem 3. In this section, we build a sketch using the

| Algorithm | Space | Threshold | Sketch time (ms) | | | | Rec time (ms) | | | | Template size (in bits) |
|--|-----------------------------------|-----------|------------------|-----------------|-----------------|---------------|-----------------|------------------|-------------------|----------------|----------------------------|
| | | | Min | Median | Max | SD | Min | Median | Max | SD | |
| SPS _{NTT} SPS _{VBB} | (\mathbb{Z}_2) ⁶⁴⁰ | 80 | 14.368 0.05 | 16.505 0.09 | 28.91 0.17 | 2.77 0.03 | 13.65 0.38 | 33.05 10.28 | 56.21 28.389 | 5.68 7.29 | 880 971 |
| SPS _{NTT} SPS _{VBB} | (\mathbb{Z}_4) ⁶⁴⁰ | 240 | 114.15 2.75 | 124.24 2.87 | 182.18 4.20 | 9.37 0.22 | 153.16 0.02 | 379.21 11.92 | 455.49 298.52 | 32.56 28.31 | 2,640 2,970 |
| SPS _{NTT} SPS _{VBB} | (\mathbb{Z}_8) ⁶⁴⁰ | 560 | 604.12 22.45 | 651.26 23.36 | 747.53 35.06 | 27.07 1.72 | 619.29 36.10 | 3,389.3 45.07 | 3,802.23 64.28 | 296.95 5.84 | 6,160 6,781 |

Table 3.2 – Experimental results for Sketch and Rec implementations over the parameters $n = 640$, $b \in \{2, 4, 8\}$, and $t = (b - 1)n/8$. Timings have been averaged over 100 iterations.

building blocks of VBB. As before, λ is a security parameter, and the integers n and t are system parameters. SPS_{VBB} can be instantiated via a tuple of probabilistic polynomial time (PPT) algorithms (ParamGen, Sketch, Rec) as described in Algorithm 3, Algorithm 4 and Algorithm 6, respectively.

The ParamGen algorithm (Algorithm 3) takes as input a security parameter λ , a positive integer n , the threshold parameter t , and a positive integer b . It outputs a family $\mathcal{G}_{n,t} = \{(\mathbb{G}, \mathbb{B}, \mathcal{B}_n)\}$. In this family, $\mathbb{G} = \mathbb{Z}_q^*$ is parametrized by prime numbers q such that \mathbb{Z}_q^* has λ -bit security with respect to DLP. Once $\mathbb{G} = \mathbb{Z}_q^*$ is fixed, \mathbb{B} is chosen as the set of all prime numbers g_i satisfying

$$\max_I \prod_{i \in I} g_i^b < \frac{q}{2} < (1 + o(1)) \max_i (g_i)^t \text{ with } |I| = \lfloor t/b \rfloor, \text{ and } I \subset \{0, \dots, n\}. \quad (3.26)$$

For a fixed \mathbb{G} and \mathbb{B} , $\mathcal{B}_n = [g_i]_{i=1}^n$ can be chosen as any ordered sequence of n pairwise distinct elements g_i in \mathbb{B} .

Algorithm 3 : ParamGen

Input : λ, n, t, b
1 Return $\{(\mathbb{G} = \mathbb{Z}_q^*,$
2 $\mathbb{B} = \text{The set of all primes } g_i \text{ satisfying (3.26),}$
3 $\mathcal{B}_n = \{[g_i]_{i=1}^n \mid g_i \in \mathbb{B} \text{ and } g_i \neq g_j \forall i \neq j\})$

The Sketch algorithm (Algorithm 4) takes as input a security parameter λ , a vector size n , a threshold t , an integer b , and a vector $x = (x_1, \dots, x_n) \in A_b^n$, where $A_b = \{0, 1, \dots, b - 1\}$. It outputs a pair (R, s) with R a parameter set and s the encoded version of x (*i.e.*, the sketch of x).

Algorithm 4 : Sketch

Input : $\lambda, n, t, b, x \in A_b^n$
Output : $R = (\mathbb{G} = \mathbb{Z}_q^*, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n), s$
1 $(\mathbb{G}, \mathbb{B}, \mathcal{B}_n) \leftarrow \text{ParamGen}(\lambda, n, t, b)$
2 $s \leftarrow \prod_{i=1}^n g_i^{x_i} \text{ mod } q$
3 Return R, s

The CFactor algorithm (Algorithm 5) takes as inputs a number x and a fixed list of primes $[g_1, \dots, g_n]$. It outputs \perp if x is composite with factors that are not in the list of primes and the factorization of x otherwise.

The Rec algorithm (Algorithm 6) receives as input R, s , and $y = (y_1, \dots, y_n) \in A_b^n$. It outputs x if $d(x, y) \leq t$ and fails (\perp) otherwise.

Algorithm 5 : *CFactor*

Input : $n, [g_1, \dots, g_n], x$
Output : S

- 1 $S = []$;
- 2 **for** $i=1, \dots, n$ **do**
- 3 **while** $g_i | x$ **do**
- 4 Append g_i to S ;
- 5 $x \leftarrow x/g_i$;
- 6 **if** $x = 1$ **then**
- 7 **Return** S ;
- 8 **else**
- 9 **Return** \perp ;

Algorithm 6 : *Rec*

Input : $R = (\mathbb{G} = \mathbb{Z}_q^*, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n), s, y$
Output : x if $d(x, y) \leq t$ and \perp otherwise.

- 1 Compute $s' = \prod_{i=1}^n g_i^{-y_i} \bmod q$;
- 2 Compute $\Delta = s \cdot s' \bmod q$;
- 3 Compute the continued fraction representation of Δ/q ,
- 4 with convergent C ;
- 5 **forall** $h/k \in C$ **do**
- 6 $F \leftarrow CFactor(n, [g_1, \dots, g_n], k)$;
- 7 $G \leftarrow CFactor(n, [g_1, \dots, g_n], k \cdot \Delta \bmod q)$;
- 8 **if** $F \neq \perp$ **and** $G \neq \perp$ **then**
- 9 Let $m = \{0\}^n$ be the zero vector;
- 10 **for** $i = 1, \dots, n$ **do**
- 11 **if** $g_i \in F$ **then**
- 12 Set m_i to minus the number of repetitions of g_i ;
- 13 **else if** $g_i \in G$ **then**
- 14 Set m_i to the number of repetitions of g_i ;
- 15 **Return** $y + m$;
- 16 **Return** \perp ;

Instantiation of SPS using NTT

NTT [71] is a cryptographic primitive to create noise-tolerant templates, where the security relies on the hardness of DLP and the Knapsack Problem (KP). SPS_{NTT} is a tuple of polynomial time algorithms (ParamGen, Sketch, Rec). The parameter generation algorithm ParamGen (Algorithm 7) takes as input a security parameter λ , a positive integer n , the threshold parameter t , and a positive integer b . It outputs a family $\mathcal{G}_{n,t} = \{(\mathbb{G}, \mathbb{B}, \mathcal{B}_n)\}$. In this family, $\mathbb{G} = \mathbb{F}_{q^2} = \mathbb{F}_q[\sigma]/\langle \sigma^2 - c \rangle$ is a finite field of size $q^2 = p^{2t}$ for some irreducible polynomial $(\sigma^2 - c)$ over \mathbb{F}_q and that \mathbb{F}_q is a finite field of size p^t parametrized by a prime number p . \mathbb{G} must have λ -bit security with respect to DLP. Once $\mathbb{G} = \mathbb{F}_q$ is fixed, \mathbb{B} is chosen as the set of all elements of the form $g = (\alpha + \sigma)/(\alpha - \sigma)$ for $\alpha \in \mathbb{F}_p$. For a fixed \mathbb{G} and \mathbb{B} , $\mathcal{B}_n = [g_i]_{i=1}^n$ can be chosen as any ordered sequence of n pairwise distinct elements g_i in \mathbb{B} .

The Sketch algorithm (Algorithm 8) takes as inputs a security parameter λ , a vector size n , a threshold t , a positive integer b and a vector $\mathbf{x} = (x_1, \dots, x_n) \in A_b^n$. It outputs (R, s) with $R = (\mathbb{G} = \mathbb{F}_{q^2}, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n)$ a parameter set and $s = \prod_{i=1}^n g_i^{x_i}$.

The Rec algorithm (Algorithm 9) takes as input $R = (\mathbb{G} = \mathbb{F}_{q^2}, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n), s, b$ a positive

Algorithm 7 : ParamGen

Input : λ, n, t, b
1 Return $\{(\mathbb{G} = \mathbb{F}_{q^2},$
2 $\mathbb{B} = \{\frac{\alpha+\sigma}{\alpha-\sigma} \in \mathbb{G} \mid \alpha \in \mathbb{F}_p\}$
3 $\mathcal{B}_n = \{[g_i]_{i=1}^n \mid g_i \in \mathbb{B} \text{ and } g_i \neq g_j \forall i \neq j\})$

Algorithm 8 : Sketch

Input : λ, n, t, x, b
Output : $(R = (\mathbb{G} = \mathbb{F}_{q^2}, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n), s)$
1 $R = (\mathbb{G}, \mathbb{B}, \mathcal{B}_n = [g_i]_{i=1}^n) \leftarrow_s \text{ParamGen}(\lambda, n, t, b);$
2 $s \leftarrow \prod_{i=1}^n g_i^{x_i} \in \mathbb{F}_{q^2};$
3 Return $(R, s);$

integer and $y = (y_1, \dots, y_n) \in A_b^n$ a vector. It then outputs x if $d(x, y) \leq t$, and \perp otherwise.

Algorithm 9 : Rec

Input : $R = (\mathbb{G} = \mathbb{F}_q[\sigma]/\langle \sigma^2 - c \rangle, \mathbb{B}, \mathcal{B}_n = [g_i = (\alpha_i + \sigma)/(\alpha_i - \sigma)]_{i=1}^n), s, b, y$
Output : x if $d(x, y) \leq t$ and \perp otherwise.
1 Compute $s' = \prod_{i=1}^n g_i^{y_i} \in \mathbb{G};$
2 $T \leftarrow s/s' \in \mathbb{G};$
3 Find $\alpha \in \mathbb{F}_q$ such that $T = ((\alpha + \sigma)/(\alpha - \sigma));$
4 For \mathbb{F}_p -variables $e_i, i = 1, \dots, t$, assign:
5 $f_0 \leftarrow \sum_{i=0}^{\lfloor t/2 \rfloor} e_{t-2i} c^i$ and $f_1 \leftarrow \sum_{i=0}^{\lfloor (t-1)/2 \rfloor} e_{t-2i-1} c^i;$
6 Use Weil restriction on $f_0 - f_1 \alpha = 0$ and solve for $e_i \in \mathbb{F}_p;$
7 Construct $P(X) = X^t + \sum_{i=1}^t (-1)^i e_i X^{t-i};$
8 Find \mathbb{F}_p -roots r_i of $P(X)$ and their multiplicity $m_i, i = 1, \dots, k.$
9 $w \leftarrow (0, \dots, 0) \in \mathbb{Z}^n;$
10 for $i \in \{1, \dots, n\}$ **do**
11 **if** $\alpha_i \in \{r_1, \dots, r_k\}$ **then**
12 $w_i \leftarrow m_j$, where $\alpha_i = r_j;$
13 **else if** $-\alpha_i \in \{r_1, \dots, r_k\}$ **then**
14 $w_i \leftarrow -m_j$, where $-\alpha_i = r_j;$
15 **else**
16 **Return** $\perp;$
17 $x \leftarrow w + y;$
18 Return $x;$

3.2.8 Robust Fuzzy Extractor based on Robust SPS

A fuzzy extractor is a cryptographic primitive designed to generate reliable and secure keys from noisy data, such as biometrics. Fuzzy extractors are composed of two main functions: Gen and Rep. The Gen function takes a noisy input w and generates a pair (s, pub) , where s is a uniformly random key and pub is public helper data. This helper data is used in the reproduction of the key. The Rep function then takes a new input, y , which is close to w , along with the public helper data, pub , and reproduces the original key, s . It is important that pub , does not divulge any substantial information about the original input w .

Using our robust SPS (see Section 3.2.5), a robust fuzzy extractor can be generically derived using [120]. Considering H a hash function indistinguishable from a random oracle, a robust fuzzy extractor (Gen, Rep) can be constructed on top of an SPS $(\text{Sketch}_{n,t}, \text{Rec})$ in the following way:

| | |
|---|--|
| $\begin{aligned} & \text{Gen}(w, \mu) \\ & \text{pub}^* \leftarrow \text{Sketch}_{n,t}(w) \\ & s \leftarrow H(w) \\ & h \leftarrow H(w, \text{pub}^*, s) \\ & \text{Return } \text{pub} = (\text{pub}^*, h), s \end{aligned}$ | $\begin{aligned} & \text{Rep}(y, \text{pub} = (\text{pub}^*, h)) \\ & w' \leftarrow \text{Rec}(y, \text{pub}^*) \\ & s' \leftarrow H(w') \\ & \text{If } H(w', \text{pub}^*, s) \neq h \text{ abort} \\ & \text{Else, return } s' \end{aligned}$ |
|---|--|

3.3 Biometric Authentication Protocol from Groups

In this section, we introduce a *zero-knowledge* authentication protocol based on SPS resistant to offline exhaustive search. More precisely, this protocol is secure against unbounded malicious clients and an honest but curious server. The protocol is given in Figure 3.2 and Figure 3.3.

Intuition Behind the Construction: The idea to protect the protocol against exhaustive offline attacks is to never provide a possible target to the client. Therefore, the protected template must remain on the server, and the authentication process relies on the ability to test whether a decomposition is correct. Specifically, the client sends its newly protected vector, and the server attempts to decompose the quotient. This way, the client can never perform an exhaustive offline search, as no target is ever provided. To ensure that interactions are indistinguishable from each other, they are randomized. The final step is to protect against the server itself. During enrollment, the client alters its secret with a nonce. Thus, even if the server reverses the protected template, it only obtains a list of possible vectors and not the actual secret, thereby ensuring the client's privacy. A final vulnerability arises if the server uses the user's newly protected vector during authentication as a target for its exhaustive search. To prevent this issue, during authentication, the client modifies its template so that the exhaustive search yields only a set of possible vectors.

The choice of k and t . To ensure *soundness*, *correctness*, and to prevent the server from trivially decomposing while allowing a genuine client to do so, the following inequalities must hold. Using the notations from Figure 3.2, we require:

1. Preventing trivial decomposability by the server requires $t < \frac{n(b-1)}{2} - k$.
2. For a genuine client to get authenticated, with e denoting the maximum distance allowed between x and y , we need $e + 2k \leq t$.
3. The *soundness* requires $\frac{n(b-1)}{4} - \frac{k}{2} > t$ to ensure that a malicious client cannot decompose.

To show that such parameters exist, here are some examples, the tuple $(b = 2, n = 128, k = 8, t = 48, e = 32)$, the tuple $(b = 3, n = 128, k = 8, t = 52, e = 36)$ and the tuple $(n = 128, b = 4, k = 8, t = 84, e = 68)$ are correct instances. The aforementioned results are direct corollaries of the following theorem.

Theorem 3.3.1. *Suppose that $v \in \{0, \dots, b-1\}$ is a random variable following a uniform distribution where $b \geq 2$ and $b \in \mathbb{Z}$. Then, the following results hold:*

1. *The expected value for v is $\frac{b-1}{2}$.*

2. With $z \in \mathbb{Z}$, the expected value of $|v - z|$ is greater than either $b/4$ with a lower bound at $z = b/2$ if b is even, or $\frac{b^2-1}{4b}$ with a lower bound at $z = (b-1)/2$ otherwise.

Proof. Let $v \in \{0, \dots, b-1\}$ be a random variable following a uniform distribution, where $b \geq 2$ and $b \in \mathbb{Z}$. Then, the expected value of v is $\mathbb{E}(v) = \sum_{i=0}^{b-1} i/b = (b-1)/2$.

Let z be an element of \mathbb{Z} , then, we consider 3 cases: $z > b-1$, $z < 0$ and $0 \leq z \leq b-1$. In the first case with $z > b-1$,

$$\mathbb{E}(|v - z|) = \sum_{i=0}^{b-1} \frac{|i - z|}{b} = z - \frac{b-1}{2} > \frac{b-1}{2}.$$

In the second case, with $z < 0$,

$$\mathbb{E}(|v - z|) = \sum_{i=0}^{b-1} \frac{|i - z|}{b} = \frac{b-1}{2} - z > \frac{b-1}{2}.$$

In the last case, with $0 \leq z \leq b-1$,

$$\begin{aligned} \mathbb{E}(|v - z|) &= \sum_{i=0}^{b-1} \frac{|i - z|}{b} = \frac{1}{b} \left(\sum_{i=0}^z (z - i) + \sum_{i=z+1}^{b-1} i - z \right) \\ &= \frac{1}{b} \left(z(2z - b + 2) + \sum_{i=0}^{b-1} i + 2 \sum_{i=0}^z i \right) \\ &= \frac{1}{b} \left(z^2 - z(b-1) \right) + \frac{b-1}{2}. \end{aligned}$$

and the result follows. ■

Corollary 3.3.1. Let $\mathbf{x} = (x_1, \dots, x_n)$ with $x_i \in \mathbb{Z}$ and $0 \leq x_i \leq b-1$ where x_i are chosen uniformly at random, then,

1. The expected weight of the vector \mathbf{x} is $\frac{n(b-1)}{2}$.
2. With \mathbf{z} a vector defined as \mathbf{x} , the expected weight of the vector $\mathbf{x} - \mathbf{z}$ is greater or equal to either $\frac{n(b^2-1)}{4b}$ if b is odd or $\frac{nb}{4}$ otherwise.

Proof. The linearity of the mathematical expectation ($\mathbb{E}(|v - z|)$) and Theorem 3.3.1 yields the result. ■

The security of our authentication scheme. Our authentication protocol ensures the *correctness*, *soundness* and *zero-knowledge* properties.

Theorem 3.3.2. The protocol depicted Figure 3.2 is sound, correct and zero-knowledge.

Proof. The server acts as the verifier and the client as the prover. For the proof, let t be the threshold.

- *Correctness:* The completeness is ensured by the *correctness* of Rec. If the client is legitimate, i.e., $d(\mathbf{x} + \rho, \mathbf{y} + \eta) \leq t$ (see the paragraph on the choice of k and t) the factorization does not fail.

Public:

- $\mathbb{G}, \mathbb{B}, n, t$.

Private:

- *Server:* $\mathcal{B}_n \subset \mathbb{B}$ and $X \in \mathbb{G}$.
- *Client:* $x, y \in A^n$.

Enrollment:

1. The client runs $\text{ParamGen}(\lambda)$ to get $A^n, \mathcal{G}_{n,t}$; and runs $\text{Sketch}_{n,t}(x)$ to get $\mathbb{G}, \mathbb{B}, \mathcal{B}_n$. It randomly draws a vector $\rho \in A^n$ such that $\sum_{i=1}^n |\rho_i| = k$ and runs the sketch on $x \oplus \rho$ to get X_ρ . The client makes \mathbb{G} and \mathbb{B} public.
2. The server stores the identity, ID , of the client, \mathcal{B}_n , and X_ρ .

Authentication:

1. The client sends its identity ID to the server.
2. The server draws randomly s in $\mathbb{Z}_{|\mathbb{G}|}^*$ and blinds the basis by computing $\mathcal{B}_n^s = [g^s : g \in \mathcal{B}_n]$. The server then sends \mathcal{B}_n^s to the client.
3. The client read its data $y \in A^n$, randomly draws a vector $\eta \in A^n$ such that $\sum_{i=1}^n |\eta_i| = k$, compute $Y_{s,\eta} = \prod_{i=1}^n g_i^{s(y_i \oplus \eta_i)} \in \mathbb{G}$ and sends it to the server.
4. The server computes $\Delta_{\eta,\rho} = \frac{(X_\rho)}{(Y_{s,\eta})^{(s^{-1})}} \in \mathbb{G}$ and runs $\text{Factor}(\Delta_{\eta,\rho})$. If the algorithm succeeds, the client is authenticated. Otherwise, the server aborts.

Figure 3.2 – SPS-based Fuzzy Authentication Scheme

- *Soundness* (Sketch of proof): Given y a vector in A^n , with the right choice of parameters, the probability that the decomposition does not fail is the probability that y is in A^n (see the paragraph on the choice of k and t). Then, the probability that a malicious prover is being accepted by the verifier is $V_t = \max_{x \in A^n} \frac{|\{y \in A^n : d(x, y) \leq t\}|}{|A^n|}$ which is negligible when n is large enough or t small enough.
- *Zero-knowledge*: The proof is based on an altered version of the protocol devoid of the first message since the verifier only wants to check if the prover knows the secret $x \in A^n$. For the proof, let A^n be the vector set, d the distance, and t the threshold.

The transcript of an iteration of the protocol is: $(\mathcal{B}_n^s, Y_{\eta,s})$.

The prover is replaced by a simulator which is given access to the secrets chosen by the server during a run of the protocol, namely s (that do not serve any purpose for the authentication). It produces a transcript of exchanges that is indistinguishable from a real back-and-forth dialogue between a genuine prover and the verifier: First, the simulator initiates the protocol. Then, it receives from the server \mathcal{B}_n^s , and uses s to recover \mathcal{B}_n . At this point, the simulator randomly draws a number m of errors from $\{1, \dots, t\}$, as well as distinct error positions $\mathcal{P} = (j_1, \dots, j_m)$ with $j_i \in \{1, \dots, n\}$. Using \mathcal{B}_n the base elements, it then computes the sketch of a simulated noisy input $Y := X \cdot (g_{j_1} \cdot g_{j_2} \cdots g_{j_m})^{-1}$ and sends $Y_s = Y^s$ to the server. Using $\Delta = X/Y$, the server computes $\text{Factor}(\Delta)$ correctly. Since the decomposition works, the checks pass successfully yielding the fake transcript. ■

Remark 3.3.1. During the run of the protocol depicted in Figure 3.2, the client received no information related to the secret during the protocol then, it did not learn anything about it.

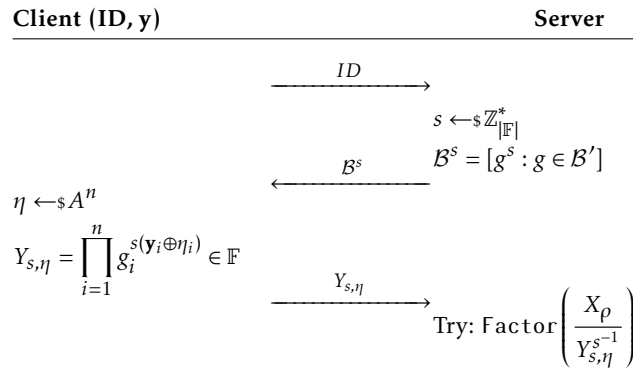


Figure 3.3 – SPS-based Authentication.

Secrecy of the enrolled reading x . The server can exhaustively search $x \oplus \rho$ by reversing the template X . The number of the possible x leading to $x \oplus \rho$ is $\binom{n}{k}$ for the binary case, yielding an information-theoretic security for x .

3.4 Remote Secure Sketch from Group

In this section, we introduce a protocol to remotely reconstruct a fuzzy input and discuss the requirements for constructing such a primitive for an SPS. Subsequently, we present a method to implement this primitive, referred to as Remote Secure Sketch (RSS). We then provide formal security arguments concerning user authentication and the secrecy of both new readings and stored template sketches.

Idea of the Construction: Secure Sketches are traditionally used locally in a single-user setup. We aim to enable the use of a Secure Sketch remotely. The main challenge is to provide the client with enough information without directly disclosing the protected template, thereby avoiding a trivial offline exhaustive search attack. The core idea is to construct a multi-party sketch. In this setup, all exchanges are blinded and randomized so that the server only knows if the client is authenticated, and the client can retrieve its secret only if they are legitimate. A schematic view of the protocol is given in Figure 3.4. In essence, the objective of the construction is to ensure the privacy of data on both ends as an oblivious functionality. As this property is difficult to obtain, we aim to ensure the full privacy of the server against a malicious client.

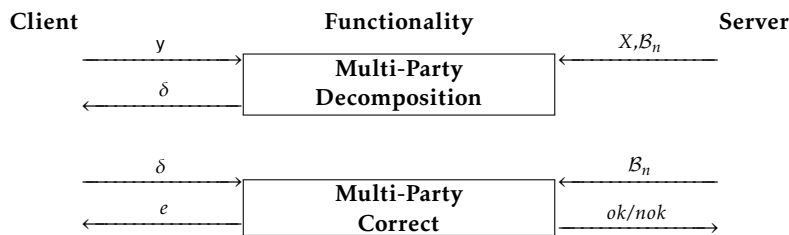


Figure 3.4 – Schematic view of SPS.

3.4.1 Requirements for Building a RSS from Subset Product Problems

We first discuss two attacks to motivate RSS. Then, an RSS protocol based on SPS is presented.

Definition 3.4.1 (Remote Secure Sketch (RSS)). *A Remote SS, or RSS, is merely a SS where the recovery function $rRec$ function ensures the following properties:*

- Given x' close from x , the client learns the corrected input $x := rRec((\mathcal{B}_n, X), x')$.
- A polynomially bounded adversary does not learn anything from the server.
- The $rRec$ algorithm is zero-knowledge meaning that the server does not learn anything useful from the client.

Constructing an $rRec$ function from subset product problems is challenging, especially when the underlying obfuscating mechanism does not fully obfuscate the desired operation. For instance, both SPS_{NTT} and SPS_{VBB} implement obfuscated distance functions, but leak the error values when the distance is below the decision threshold. However, an obfuscated distance function that leaks information only locally is useful for constructing a secure sketch. The challenge here is twofold: (i) Allowing the legitimate client to correct her fresh measurement x' without leaking to the server any information about the enrolled reading x (resistance to accumulation attacks described in Section 2.4.5); (ii) Preventing polynomially bounded adversaries from conducting offline exhaustive search attacks. This is achieved by preventing them from learning any information by interacting with the server.

3.4.2 Mechanics of our Remote Subset Product Sketch

This section describes the oblivious SPS sketch which can be instantiated using either VBB or NTT. Enrollment (storage of the sketch) is illustrated in Figure 3.6, while the remote recovery is illustrated in Figure 3.7. The main idea is to make the Rec algorithm collaborative.

RSS from SPS. We construct a Remote Subset Product Sketch (RSPS) that is resistant to both offline and accumulation attacks. It is a pair (Sketch, $rRec$) of online randomized procedures which respectively enable the generation of a helping value and the recovery of the original (enrolled) data by collaborating with a server. The online procedures are detailed in Figure 3.5. Note that during each step of the protocol, both parties expect non-trivial elements from the group. Failure to fulfill this condition results in the abortion of the protocol.

Preventing Offline Attacks. For conducting an exhaustive search, a malicious client requires two things, a way to encode an arbitrary vector and a way to check her guess using a target. In other words, she needs the elements of the basis \mathcal{B}_n along with the sketch template X . The protocol, depicted in Figure 3.5, is designed to prevent giving any target for an offline exhaustive search. Specifically, the malicious client has no access to the basis since the server blinds it with a random exponent. To get the sketch X , an adversary can run the protocol and submit a particular template, e.g., $(0, \dots, 0)$. However, this is detectable by the server, and even if undetected, the product $H = \prod_{i=1}^k h_i$ blinds the sketch. This co-factor is also here to prevent leakage of the basis when the client is malicious. Without H , if the client makes two queries with x and y such that $d(x, y) = t$, it is possible to deduce t basis elements. More precisely, let Δ_x denote the Δ obtained by querying x and let Δ_y denote the one obtained by querying y . Then, factorizing Δ_x/Δ_y leaks the g_i where $x_i \neq y_i$. Once \mathcal{B}_n is entirely leaked after $\lceil n/t \rceil$ queries, an offline exhaustive search attack

Public:

- $\mathbb{G}, \mathbb{B}, n, t$.

Private:

- Server: $\mathcal{B}_n \subset \mathbb{B}$ and $X \in \mathbb{G}$.
- Client: $x, y \in A^n$.

Sketch:

1. The client runs $\text{ParamGen}(\lambda)$ to get $A^n, \mathcal{G}_{n,t}$; and runs $\text{Sketch}_{n,t}(x)$ to get $\mathbb{G}, \mathbb{B}, \mathcal{B}_n$ and the sketch X . The client makes \mathbb{G} and \mathbb{B} public.
2. The server stores the identity, ID , of the client, \mathcal{B}_n , and X .

rRec:

1. The client sends its identity ID to the server.
2. The server draws randomly r and s in $\mathbb{Z}_{|\mathbb{G}|}^*$ and samples with replacement k basis elements $h_1, \dots, h_k \in \mathbb{B} \setminus \mathcal{B}_n$. The server blinds the basis and the template by computing $H = \prod_{i=1}^k h_i \in \mathbb{G}$, $X_r = X^r H^r \in \mathbb{G}$ and $\mathcal{B}_n^s = [g^s : g \in \mathcal{B}_n]$. The server then sends \mathcal{B}_n^s and X_r to the client.
3. To blind its data, the client draws randomly γ in $\mathbb{Z}_{|\mathbb{G}|}^*$, computes $X_{r,\gamma} = X_r^\gamma \in \mathbb{G}$ and $\prod_{i=1}^n g_i^{s \cdot r \cdot \gamma_i} \in \mathbb{G}$. $X_{r,\gamma}$ and $Y_{s,\gamma}$ are sent to the server.
4. The server computes $\Delta_\gamma = \frac{(X_{r,\gamma})^{(r^{-1})}}{(Y_{s,\gamma})^{(s^{-1})}} \in \mathbb{G}$ and sends it to the client.
5. The client computes $\Delta = \Delta_\gamma^{\gamma^{-1}} \in \mathbb{G}$ and runs $\text{Factor}(\Delta)$ to obtain an ordered sequence of integers $[\delta_k]_{k=1}^{|\mathbb{B}|}$ such that $\Delta = \prod_{k=1}^{|\mathbb{B}|} g_k^{\delta_k}$, where $g_k \in \mathbb{B}$, and $\sum_{k=1}^{|\mathbb{B}|} |\delta_k| \leq t$. The client computes $U = [g_i : \delta_i \neq 0, i \in \{0, \dots, |\mathbb{B}|\}]$ and sends it to the server.
6. The server first checks if $U \setminus \{h_1, \dots, h_k\} \subseteq \mathcal{B}_n$, and then, for $i = 1, \dots, k$, checks if $h_i \in U$. If one of the checks fails, the server aborts. The server computes $\mathcal{P} = \text{Get_Index}(U)$, which returns a vector containing the index of every element from U within \mathcal{B}_n . If the element is not in \mathcal{B}_n , its position is -1 . The server sends \mathcal{P} to the client.
7. The client retrieve x by using the fact that $x_i = y_i$ if $i \notin \mathcal{P}$ and $x_i = y_i + \delta_i$ otherwise.

Figure 3.5 – OSS scheme based on SPS

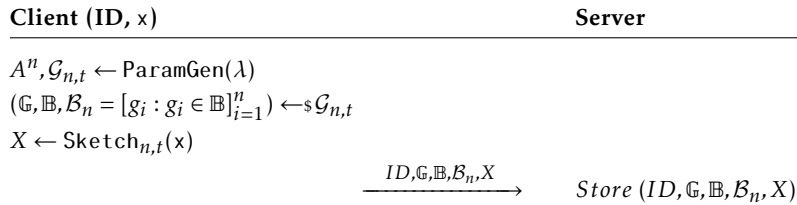
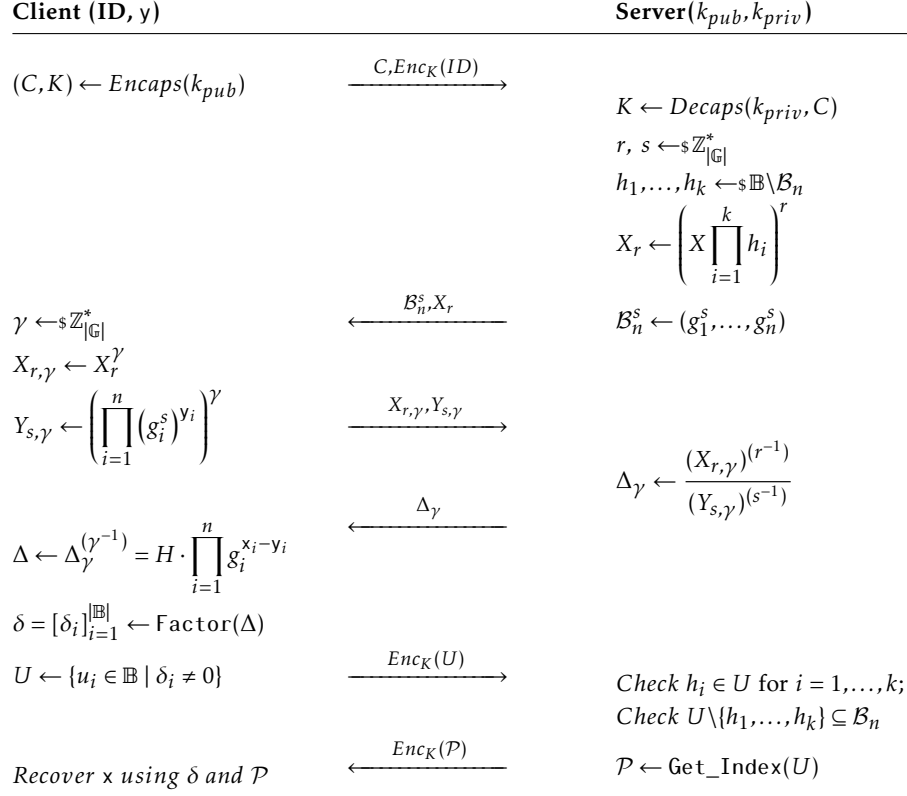


Figure 3.6 – SPS-based Sketch.

Figure 3.7 – SPS-based $r\text{Rec}$.

may be possible if X is known by the attacker. The presence of H is detailed in the *soundness* of the *zero-knowledge* proof, but, in summary, the presence of the factorization of H in the last message sent to the server ensures that the client is legitimate. If the client is not legitimate, the server does not answer the last message of the malicious client (used to check if some elements belong to the basis), thus preventing the leakage of the chosen basis elements.

The choice of k and t . The *correctness* is ensured if the primitive SPS is used with a threshold $t \leq n/2$. Furthermore, if e_{\max} denotes the maximum number of errors allowed for a genuine client, then $e_{\max} + k \leq t$ is necessary for unique t -factorization. On the other hand, to prevent active adversarial attacks (e.g. collecting Δ , Δ' ; and factorizing Δ/Δ'), it is necessary that $2k > t$. Therefore, we have $e_{\max} + k \leq n/2 < 2k$. For a fixed n , one would have $k > n/4$ and $e_{\max} < n/4$.

Preventing accumulation Attacks. The protocol described in Figure 3.7 does not leak extra information to the server. The accumulation attack is possible when the server knows the position and the value of some errors. However, thanks to the Factor algorithm, the server learns only the positions of the errors, and not the exact values of the errors at these positions. Hence, the server does not know whether $x_i > y_i$ or not for any i , as γ is used to prevent the server from factoring.

Remark 3.4.1. Let us suppose a Key Encapsulation Mechanism (KEM) made of a triple of algorithms: The key generation (KeyGen) generates the pair of private and public keys (k_{priv}, k_{pub}) ; The encapsulate algorithm (Encaps) takes the public key (k_{pub}) and produces a ciphertext (C) to be passed to the recipient and the shared secret key K for the originator; the decapsulate algorithm (Decaps) takes the recipient

private key (k_{priv}) and the ciphertext (C) and produces the shared secret key (K). In the recovery Figure 3.7, the exchanged messages should be encrypted using a shared secret key K . For this purpose, before sending her ID, the client uses a KEM to encapsulate and send a symmetric key to the server using its public key. This key K serves as input for symmetric encryption (denoted Enc) to encrypt the first and two last messages. The encryptions and blindings make the transcripts indistinguishable, hence preventing the tracking of users.

3.4.3 Authentication and Secrecy of Fresh Readings

In the rRec procedure, the client needs to authenticate her fresh reading to the server before the server returns any information allowing the recovery of the enrolled reading. We show that the function ensures the properties of a zero-knowledge proof of knowledge.

Theorem 3.4.1. *The rRec protocol depicted Figure 3.5 is complete.*

Proof. The server is acting as the verifier (V) and the client as the prover (P). The proof is based on an altered version of the protocol devoid of the last exchanged message since the verifier only wants to check if the prover knows the secret $x \in A^n$. For the proof, let t be the threshold.

The completeness is ensured by the *correctness* of Rec. If the client is legitimate, i.e. $d(x, y) \leq t$, the factorization of δ is correct and she can correct his data. ■

Theorem 3.4.2. *The interactive proof of the rRec protocol depicted Figure 3.5 is sound.*

Sketch of a proof. The server is acting as the verifier (V) and the client as the prover (P). The proof is based on an altered version of the protocol devoid of the last exchanged message since the verifier only wants to check if the prover knows the secret $x \in A^n$. For the proof, let A^n be the vector set, d the distance, and t the threshold.

If P tries to run the protocol without the knowledge of x , she has two strategies:

1. Picking y a random vector in \mathcal{B} . She then succeeds the protocol with probability $\frac{|\{y \in A^n | d(y, x) \leq t\}|}{|A|^n}$.
2. Trying to present to the server a list U' such that all the chosen $(h_i)_{i=1, \dots, k}$ are in U' . Let $\overline{\mathcal{B}}_n = \mathbb{B} \setminus \mathcal{B}_n$ for a given base \mathcal{B}_n (chosen by the client during enrollment). The success probability is less than $\frac{k^k}{|\overline{\mathcal{B}}_n|^k}$.

If needed, we can add a ZKP on the reconstructed x with the storing of $H(x)$ on the server, and the proof follows. ■

Theorem 3.4.3. *The rRec protocol depicted Figure 3.5 achieves computational zero knowledge².*

Proof. The server is acting as the verifier (V) and the client as the prover (P). The proof is based on an altered version of the protocol devoid of the last exchanged message since the verifier only wants to check if the prover knows the secret $x \in A^n$. For the proof, let A^n be the vector set, d the distance, and t the threshold.

The transcript of an iteration of the protocol is:

$$(ct, Enc_K(ID, K), X_r, \mathcal{B}_n^S, X_{r,\gamma}, Y_{s,\gamma}, \Delta_\gamma, Enc_K(U), Enc_K(\mathcal{P}))$$

2. Computational zero knowledge means that no there are no efficient algorithm to distinguish the fake transcript distributions and the genuine transcript.

The prover is replaced by a simulator which is given access to the secrets chosen by the server during a run of the protocol, namely r and s (that do not serve any purpose for the authentication). It produces a transcript of exchanges that is indistinguishable from a real back-and-forth dialogue between a genuine prover and the verifier: First, the simulator initiates the protocol. Then, it receives from the server \mathcal{B}_n^s and X_r , and uses s and r to recover \mathcal{B}_n and X . At this point, the simulator randomly draws a number m of errors from $\{1, \dots, t\}$, as well as distinct error positions $\mathcal{P} = (j_1, \dots, j_m)$ with $j_i \in \{1, \dots, n\}$. Using \mathcal{B}_n the base elements, it then computes the sketch of a simulated noisy input $Y := X \cdot (g_{j_1} \cdot g_{j_2} \cdots g_{j_m})^{-1}$. Next, the simulator chooses $\gamma \leftarrow \mathbb{Z}_{|G|}^\times$, and then computes $Y_{s,\gamma} := Y^{s\gamma}$ and $X_{r,\gamma} := (X_r)^\gamma$ which are both sent to the server. After, the simulator receives Δ_γ , and computes U with $\text{Factor}(\Delta_\gamma^{\gamma^{-1}})$. Since U contains all the elements of the product H along with m elements of \mathcal{B}_n , the checks pass successfully. The end of the protocol produces the end of the fake transcript. ■

Remark 3.4.2. *The statistical difference is present on Y . For a genuine client, the exponent values of the g_i that compose Y are between 0 and $q - 1$. In the case of the simulator, the exponent values are between -1 and q . In order to differentiate between these two distributions, it is necessary to reverse the Y template or to solve the DSPP. Consequently, no polynomial algorithm exists for distinguishing between the two distributions.*

3.4.4 Secrecy of the Template Sketch

Preventing offline exhaustive search attacks requires maintaining the secrecy of $X = \prod_{i=1}^n g_i^{x_i}$ throughout the interactions. We show that the secret is preserved while the basis \mathcal{B}_n is exposed to an adversary.

Two formal arguments are provided. First, we show that, when interacting with semi-honest adversaries, the server can be replaced by a simulator while yielding transcripts that are indistinguishable from real transcripts.

Theorem 3.4.4. *While interacting with adversaries that do not deviate from the rRec protocol and who knows \mathcal{B}_n , the server preserves the secrecy of the template sketch X .*

Proof. There exists a simulator for a client that is indistinguishable from a real view in a real interaction. The client produces a valid transcript without the knowledge of x and X while interacting with the simulator. The client and the simulator run the protocol as usual until the simulator gets $Y_{\gamma,s}$. It then recovers Y as the randomizers s and γ are known. The simulator draws randomly $\varepsilon \in A^n$ such that $\sum_{i=1}^n \varepsilon_i \leq t$ and go back in time to set $X = Y \prod_{i=1}^n g_i^{\varepsilon_i}$. The execution of the protocol now provides an indistinguishable trace. ■

Remark 3.4.3. *We believe that the basis remains secret while interacting with the client. This provides more security as the knowledge of both the basis and the template sketch is required to enable offline exhaustive search attacks.*

Resistance to offline search attacks under the compromise of X and g_i . Consider a malicious client having the knowledge of X as well as the basis elements g_i up to their order. Even if the attacker is computationally unbounded, offline attacks are not possible. In other words, even if the attacker is all-powerful, she may not retrieve x without interacting with the server.

Public:

- k_{pub} an asymmetric public key and Ext an extractor.

Private:

- Server: HD the helping data and μ a salt.

Gen:

1. The client reads his data $x = (x_1, \dots, x_n)$ and applies Sketch on x to get HD the helping data. The client draws randomly μ its salt and generates its secret $Ext(\mu, x)$.
2. The server stores the identity ID of the client, the helping data HD and the salt μ .

Reproduce:

1. The server and the client run the RSS protocol.
2. If the protocol does not abort, the server sends μ to the client.
3. Then she reproduces her secret $Ext(\mu, x)$.

Figure 3.8 – Generic construction of an oblivious fuzzy extractor (RFE) on top of an oblivious secure sketch scheme (RSS).

Theorem 3.4.5. *A computationally unbounded adversary who knows X , along with the $g_i \in \mathcal{B}_n$ up to their order, cannot retrieve the enrolled reading x among the exponential number of solutions.*

Proof. Suppose that the attacker can reverse the template X by assuming an order of the g_i . She then gets the list of components x_1, \dots, x_n up to their order. Assuming that the readings x are uniformly distributed in A^n , the number of possible orders leading to distinct solutions is

$$\frac{n!}{((n/b)!)^b} \approx b^{n-o(\log n)}$$

for a constant $b = |A|$ and a large n . ■

3.5 Secure Remote Fuzzy Extractors from Remote Secure Sketches

A straightforward way to build a fuzzy extractor is to use an extraction function such as a universal hash function and a secure sketch (see Section 3.2.8 and [120]). Thus, it is natural to construct a secure remote fuzzy extractor by using a remote secure sketch and a universal hash function. Figure 3.8 provides a generic method to build a remote fuzzy extractor.

Definition 3.5.1 (Remote Fuzzy Extractor). *An Remote Fuzzy Extractor (RFE) is a Fuzzy Extractor (FE) where the Reproduce function, denoted rRep ensure the following properties:*

- Given x' close from x , the client learns the correct key $s := rRep(x', HD)$.
- A polynomially bounded adversary does not learn anything from the server.
- The rRep algorithm does not leak any information to the server.

The generic construction of a RFE from a RSS is given in Figure 3.8.

Remark 3.5.1. *As RFE schemes are based on RSS schemes, they are resistant to both offline and accumulation attacks under the same hypothesis. If we do not want the client to learn anything about the salt, the extraction function $Ext(\cdot)$ can be replaced by an Oblivious Pseudo Random Function [28].*

3.6 Conclusion and Future Work

In this chapter, we present a method to construct groups with a unique factorization property. It would be interesting to investigate how to construct such cryptographic groups other than through NTT and VBB. A construction proof might be useful to seek new attacks on such constructions. Moreover, the unique factorization property could be useful for future work in related areas *e.g.*, for secret sharing schemes. However, as the uniqueness of the decomposition relies on the Gaussian Heuristic, more investigations are needed to quantify the probability that the hypothesis fails and completely assess the security of our construction.

Based on this construction, we propose a generic method for constructing secure sketches from groups, which we refer to as SPS. The SPS algorithm operates on integer vectors for the L_1 distance, thereby extending previously known sketch constructions on binary vectors with the Hamming distance. We demonstrated that the proposed construction is reusable, irreversible, and robust. We provided concrete instantiations of the construction based on previously known primitives, namely, NTT and VBB. Based on the state-of-the-art solvers for discrete logarithm, knapsack, and short vector problems, we provided bit-level security estimates for the construction. Future work involves a theoretical description of the computational cost of SPS. Another direction for further investigation is the enhancement of SPS to handle private intersections. In this manner, our construction may be used as the foundation for fuzzy vault schemes. It would also be interesting to challenge the security of SPS and find new attack strategies with better complexities.

To prevent offline exhaustive search attacks, we propose a *zero-knowledge* protocol to prove the knowledge of a fuzzy input as well as the notion of remote secure sketches to remotely reconstruct fuzzy data. As SPS sketches are based on groups, standard blinding techniques have been applied to derive an instantiation of the construction. Future research will focus on identifying alternative methods for generating both remote secure sketches and remote fuzzy extractors relying on theoretical information security. Some tools have been developed in the literature that could assist in achieving this objective. Such techniques include those of secure MultiParty Calculation (MPC [55, 138, 90]) and Oblivious Transfer (OT [141, 121]). During the development of our protocol, we were unable to simultaneously achieve total resistance to offline attacks (*i.e.*, information theoretical security) and robustness. It is reasonable to believe that these two properties cannot be achieved simultaneously by our method. Indeed, to resist offline attacks against an attacker who is not computationally bound, it is necessary to ensure that the attacker is not provided with any information about the secret protected by a computational hypothesis. Concurrently, a sufficient quantity of information dependent on the secret and the information possessed by the server must be employed to ascertain the identity of the server for robustness. It is evident that, given the interactive nature of our construction, the client must abort the connection before responding to the server, thereby realizing that the server is malicious. The issue is that the secret is required to ascertain whether the server is malicious. The secret is only obtained at the end of the interaction. Consequently, robustness cannot be achieved. The aforementioned reasoning appears to apply to all secure remote secure sketches and remote fuzzy extractors in general. Further research is required to ascertain whether a client can identify a malicious server without the use of a second factor or an additional party, such as an authority.

Conclusion and Future Work

Outline of the current chapter

| | |
|-----------------------------|-----|
| Conclusion | 101 |
| Future Work | 102 |

Conclusion

In this manuscript, we have conducted an in-depth study of biometrics through three major axes, each addressing distinct challenges.

The first axis, **Security Assessment**, focused on the comprehensive evaluation of biometric systems and transformation schemes. Given the diverse nature of biometric technologies and the limited commonalities among existing systems, this axis posed a significant challenge. To tackle this, we analyzed biometric systems at the most fundamental level: the template space. Our investigation revealed that near-collisions are a key factor leading to the loss of precision in biometric systems. By studying near-collisions, we elucidated the mechanisms governing their occurrence and developed an attack model capable of simultaneously impersonating multiple individuals. We also devised a novel scoring metric to help configure biometric systems in a way that minimizes the impact of near-collisions. Furthermore, we explored brute-force attacks, which are universally applicable across different biometric systems. Our characterization of potential attack strategies provides a foundational understanding of biometric system security and serves as a reference for assessing more specific attacks. This work is expected to enhance the comparability of biometric systems in terms of security, thereby increasing competition among research teams and contributing to an overall improvement in biometric system security.

The second axis, **Identifying and Refining Low-Level Primitives**, was dedicated to discovering and improving fundamental building blocks for use in biometric protocols. This axis was particularly challenging due to its exploratory nature and the mathematical complexity involved, with the possibility of unsuccessful outcomes. Despite these challenges, we developed a cryptographic primitive that evaluates the distance between two vectors in a homomorphic manner. This structure also allows for data correction to match reference data if the distance is below a specified threshold. The main advantage of this construction lies in its computational hardness, making it difficult to recover reference data or link two enrollments. Although primarily designed for biometric applications, this cryptographic primitive is also applicable to other domains dealing with fuzzy data, such as secret sharing schemes, thus broadening its potential impact.

The final axis, **Describing Recognition Protocols**, aimed to design biometric authentication

protocols that are resistant to attacks while preserving user privacy and minimizing computational and communication costs. To distinguish our work from existing constructions, we focused on achieving a property previously unattainable in scenarios where the client relies solely on biometric data: resistance to offline exhaustive search attacks. Our first protocol addresses this challenge by utilizing a zero-knowledge proof mechanism, ensuring that a computationally unbounded malicious client cannot extract useful information during the authentication process. Additionally, we introduced a second protocol that facilitates both authentication and the reconstruction of the original secret. This protocol is designed for scenarios where the adversary is computationally bounded, providing a balance between security and practical feasibility.

In summary, this research has made significant contributions to the field of biometric security by providing a thorough analysis of biometric vulnerabilities, developing innovative cryptographic primitives, and designing advanced authentication protocols. These advancements offer practical solutions to enhance the security and efficiency of biometric systems, laying a solid foundation for future research and development in this field.

Future Work

In this section, future work related to the challenges identified in biometric systems is presented, focusing on the analysis and formalization of biometric accuracy metrics and third-party removal of unlinkability.

Analyze and Formalize Biometric Accuracy Metrics: As a result of our research, it has become evident that biometric accuracy metrics such as FMR, FNMR, FRR, and FAR are either not well-defined or not correctly utilized. Although some work [97, 102] in the literature attempts to formalize and provide a statistical analysis of FMR, FNMR, and EER, these contributions are not widely known within the field. We believe there is significant potential to improve these results and to clarify them in the simplest way possible. By doing so, we aim to promote the adoption of robust methodologies in the field, similar to the impact of Shoup’s tutorial on game-based proofs [126]. Indeed, for a significant proportion of authors FMR is identical to FAR and FNMR to FRR despite their notable differences (see Section 1.3.2). Furthermore, authors in the literature do not provide the methodology used to calculate these metrics, making comparisons challenging or even uninterpretable. For future research, it would be beneficial to develop a rigorous and systematic methodology for defining, calculating, and interpreting both those accuracy metrics. In the following, we explain how we project to do so by taking the FMR as an example, but it is trivially generalized to other metrics. To achieve this, it is necessary to conceptualize the FMR as an unknown probabilities that need to be estimated. This requires the establishment of an estimator of this probability. To do so, statistical methods must be used for estimating FMR involving confidence intervals to quantify the uncertainty associated with the estimator. More precisely, an empirical estimate of the false match rate denoted by $\overline{\text{FMR}}$ can be calculated by dividing the number of false matches by the number of matches tested for a given system. More precisely, according to the Face Recognition Technology Evaluation (FRTE) 1:1 Verification [14], given a vector of N imposter scores v and T a threshold, an estimation of the FMR is

$$\overline{\text{FMR}} = \frac{1}{N} \sum_{i=1}^N S(T - v_i) = \frac{\text{Number of false matches}}{\text{Total number of comparisons}} \quad (3.27)$$

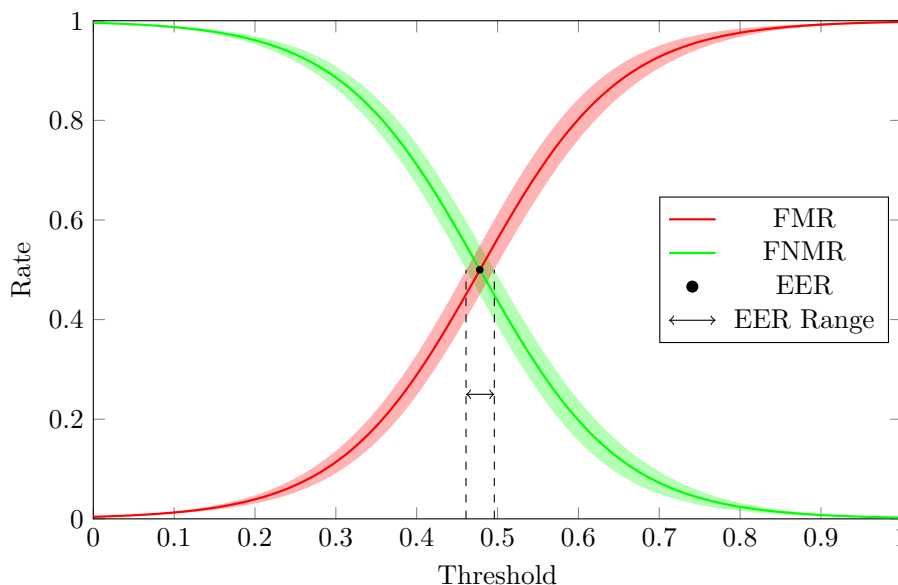


Figure 3.9 – EER interval from the uncertainty of the FMR and FNMR.

with $S(\cdot)$ the unit step function, and $S(0)$ taken to be 1. $S(T - v_i)$ can only take two values *i.e.*, 0 or 1 meaning that either the corresponding pair of users match or not. By convention, if we consider that matching is a success then, $S(T - v_i) \sim X_i$ is a random variable that follows a Bernoulli distribution of probability FMR. The above equation can be rewritten as

$$\overline{\text{FMR}} = \frac{1}{N} \sum_{i=1}^N X_i. \tag{3.28}$$

Then, using the Central Limit Theorem [65] which states that under reasonable assumptions on the unknown distribution of the data, as long as the sample size is large, the distribution of the sample mean tends almost surely to be normally distributed, our estimator is the mean of a random variable which is normally distributed. As the normal distribution is well studied, we can compute a Confidence Interval [59] (CI) on $\overline{\text{FMR}}$. The CI at 95% on $\overline{\text{FMR}}$ computed from n samples is given by

$$\text{FMR} \approx \overline{\text{FMR}} \pm 2 \sqrt{\frac{\overline{\text{FMR}} \times (1 - \overline{\text{FMR}})}{n - 1}}. \tag{3.29}$$

Then, with the same methodology on the FNMR, we are able to provide an interval for the threshold yielding the Equal Error Rate (EER) as depicted in Figure 3.9. This statistics analysis would provide enough tools to ascertain the real security of a biometric system. The anticipated outcome of this research is the establishment of standardized methodologies for the calculation of biometric metrics especially the EER, with the ultimate objective of developing systems that can be readily compared.

Modeling and analyzing the Zombie Attack: A zombie computer [93], often simply referred to as a "zombie," is a compromised computer that has been infiltrated by a malicious actor, typically via malware. Once compromised, the computer is remotely controlled without the knowledge or consent of the user. We then propose the following attack model. An attacker tries to exhaustively search passwords within a database. Each successful attempt infects a user, turning their computer into a zombie working for the attacker. The idea is to study how long it would take for such an

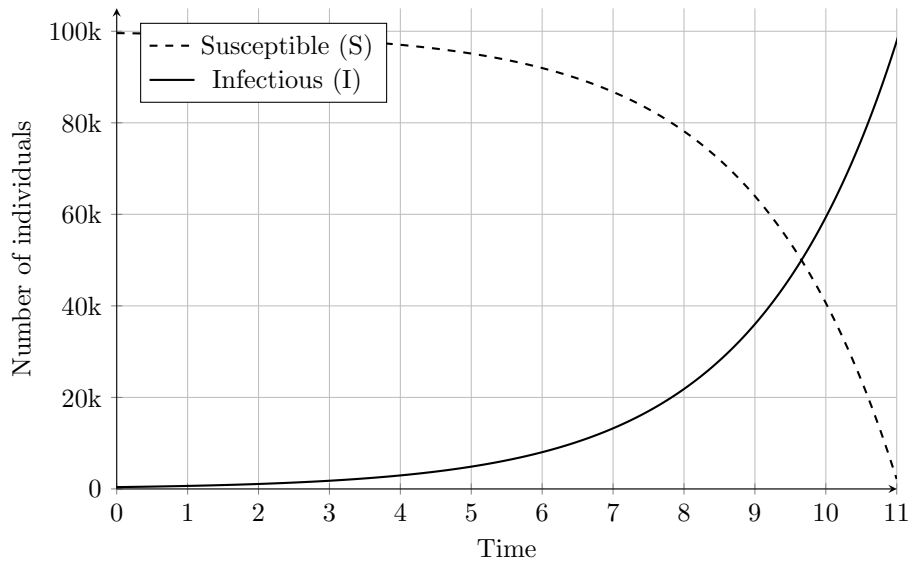


Figure 3.10 – Illustration of a Susceptible-Infectious (SI) model for the zombie attack.

attacker to corrupt all the users in a given database. This attack directly applies to all Windows users (more than 70% of worldwide users [10]) as their Microsoft account password is used to unlock their devices. Our first goal is to build a simulator to mimic real-world environments where these attacks may occur, considering varying levels of password security enhancements. We would use statistical models such as SIR (Susceptible-Infectious-Recovered), SIS (Susceptible-Infectious-Susceptible), and SEIS (Susceptible-Exposed-Infectious-Susceptible) from the outset to simulate the dynamics of cyberattack propagation. The choice of the SIR, SIS, and SEIS models is motivated by their proven effectiveness in capturing the spread dynamics of infectious diseases, which can be analogous to the propagation of malware in a network as shown in [36, 64, 89]. In the SIR model, the population is divided into three compartments: Susceptible (S), who are vulnerable to infection; Infectious (I), who have been infected and can spread the malware; and Recovered (R), who have recovered from the infection and are assumed to be immune. The transition from Susceptible to Infectious occurs at a rate proportional to the number of contacts between susceptible and infectious individuals, while the transition from Infectious to Recovered happens at a certain recovery rate. The basic reproduction number, R_0 , indicates how many new infections one infected individual would cause in a completely susceptible population. In order to accurately estimate the parameters of the SIR model, it is necessary to conduct a thorough review of the existing literature on virus propagation [45], and to perform direct simulations by launching an attack on a local database such as "rockyou" [35]. An illustration of the SI model is provided in Figure 3.10. Another potential enhancement is to simulate the impact of various cybersecurity measures, such as multi-factor authentication (by giving multiple passwords per user that needs to be found before getting infected), to evaluate their effectiveness in slowing down or preventing the spread of zombie attacks. The objective of this work is to develop a comprehensive and reliable tool for the analysis and mitigation of risks associated with zombie attacks. This tool would significantly contribute to the enhancement of the security of user devices and networks.

Résumé par chapitre

Dans cette partie du manuscrit, nous résumons les contributions des chapitres précédents. Chaque section suit le même schéma. Tout d'abord, nous présentons un résumé du chapitre en question. Ensuite, nous remettons notre travail dans son contexte. Enfin, nous indiquons les contributions détaillées.

Chapitre 2

Résumé : Dans ce chapitre, nous examinons la sécurité des données biométriques face à différentes attaques. Les données biométriques, après transformation, se comportent de manière similaire aux mots de passe. Par conséquent, nous les analysons de la même manière. L'analyse implique donc trois aspects essentiels que nous appliquons également ici.

Le premier aspect est la recherche de collisions. Une différence majeure entre les données biométriques et les mots de passe est la variabilité intrinsèque des données biométriques, même pour une même personne. Par conséquent, dans le cas des données biométriques, le système doit tolérer les variations et les correspondances exactes ne sont pas pertinentes, ce qui est la norme pour les mots de passe. Ainsi, au lieu de collisions exactes, nous considérons des quasi-collisions. Nous les étudions et les caractérisons car elles posent un problème à la fois pour la performance et la sécurité des systèmes biométriques. Un grand nombre de quasi-collisions dégrade considérablement la précision du système. De plus, nous montrons comment un attaquant peut exploiter les quasi-collisions pour créer des gabarits capables d'usurper l'identité de plusieurs utilisateurs. Pour atténuer ce problème, nous déterminons leurs probabilités d'occurrence et dérivons une méthode de score pour aider à la paramétrisation des systèmes biométriques.

Le deuxième aspect concerne les attaques ciblées par recherche exhaustive, qui visent un utilisateur spécifique. Bien que ce travail ait été initié dans la littérature, nous l'étendons pour prendre en compte différents scénarios d'attaques. Nous caractérisons divers scénarios de fuite de données, tels que la fuite de distance ou la fuite des positions des erreurs. Chaque scénario conduit à des attaques génériques dont nous fournissons les analyses de complexité. Cela permet de concevoir des systèmes biométriques avec une compréhension claire de la complexité réelle des attaques par recherche exhaustive en fonction des fuites d'informations du matcheur.

Le dernier aspect concerne les attaques non ciblées par recherche exhaustive, qui ne visent aucun utilisateur en particulier. Cela est crucial car dans une base de données très peuplée, même si la probabilité d'usurper l'identité d'un individu spécifique reste faible, la probabilité d'usurper l'identité de quelqu'un augmente significativement. Nous caractérisons la probabilité qu'un individu dans la base de données voit son identité usurpée. Cette caractérisation nous permet de dériver un score de sécurité contre de telles attaques et de fournir des recommandations pour limiter leur impact.

Contexte : Les systèmes biométriques peuvent être divisés en deux catégories principales. La première utilise des données en clair pour faciliter la prise de décision. Ces types de systèmes ne sont pas conçus pour préserver la vie privée de l'utilisateur, mais ils sont (très) précis. Dans ce chapitre, nous nous concentrons sur le deuxième type de système, dans lequel les données biométriques ne sont pas en texte clair. Bien que ces systèmes soient moins précis que leur homologue, ils garantissent néanmoins la protection de la vie privée de l'utilisateur. Pour se faire, ces algorithmes utilisent une fonction de transformation qui prend les données biométriques en entrée, avec

ou sans informations supplémentaires, et produit des données biométriques protégées appelées gabarits. Ces transformations doivent respecter les spécifications ISO [1, 2]. Ces spécifications comprennent quatre propriétés essentielles. *Irréversibilité* : Étant donné un gabarit, il doit être difficile de récupérer les données biométriques. Cette propriété peut être généralisée. Ainsi, on veut qu'étant donné un gabarit, il doit être difficile de trouver une donnée biométrique qui, par la même transformation, donne le même gabarit ou un gabarit proche vis-à-vis du seuil de décision du système. *Impossibilité de liaison* : Étant donné deux gabarits différents, il doit être difficile de déterminer s'ils proviennent ou non du même utilisateur. *Révocabilité* : Il s'agit de la possibilité de révoquer, d'annuler ou de désactiver un gabarit compromis et d'en générer un nouveau sans affecter la sécurité du système. La dernière propriété est la *préservation des performances*, qui stipule qu'une transformation ne doit pas entraîner une perte significative de précision. L'objectif de cette transformation est d'utiliser les données biométriques comme mots de passe tout en préservant leur confidentialité (voir Section 1.3.1 pour plus de détails). De manière classique, la protection des mots de passe est assurée par l'utilisation de fonctions de hachage qui respectent deux propriétés majeures : l'*irréversibilité* et la *résistance aux collisions*. Ces deux propriétés sont identiques à celles de l'*irréversibilité* pour les gabarits, sauf que pour les hachés, l'attaquant recherche des correspondances parfaites. En outre, les hachés peuvent facilement être révoqués en modifiant le sel ou le mot de passe. Compte tenu de leurs propriétés et applications analogues, il est raisonnable de considérer *de facto* que les gabarits s'apparentent à des hachés et qu'ils doivent être étudiés au même titre. Dans la littérature, il existe trois approches principales de l'étude des fonctions de hachage. La première consiste à étudier l'occurrence des collisions. La deuxième, examine la complexité d'une attaque contre un utilisateur spécifique, c'est-à-dire les attaques ciblées. La troisième examine la sécurité du système lorsque tous les utilisateurs sont attaqués simultanément, c'est-à-dire les attaques non ciblées. Dans ce chapitre, nous appliquons cette méthodologie aux gabarits biométriques afin d'estimer et de borner leur sécurité en cas d'attaques par quasi-collision, ciblées et non ciblées.

Contributions détaillées : Les contributions présentées dans la liste suivante donnent un aperçu des différentes facettes abordées dans notre travail, offrant une compréhension globale des résultats obtenus.

- Section 2.2 : Notre première contribution est un algorithme de partitionnement efficace qui accélère les attaques visant à générer des gabarits maîtres. Ceux-ci sont des gabarits qui permettent d'usurper l'identité de différents individus simultanément. Les études numériques que nous avons menées sur la mise en œuvre de notre méthode démontrent une réduction du temps de calcul dans certains contextes. De plus, nous mettons en évidence un lien entre le problème de la chaîne la plus proche (closest string problem) avec un nombre arbitraire de mots et les gabarits maîtres. Ainsi, nous fournissons une solution à ces deux problèmes avec la même méthodologie basée sur le recuit simulé (SANN). En outre, nous déterminons une limite à la taille d'une base de données en fonction de la dimension de l'espace des gabarits et du seuil de décision, ce qui permet d'éviter les quasi-collisions avec une probabilité élevée. Nous introduisons la notion de quasi-collision faible et de quasi-collision forte, ce qui nous permet de fournir une analyse théorique de la sécurité des schémas de transformation biométrique. Les bornes sur la probabilité d'une quasi-collision mettent en évidence les limites théoriques de la précision d'un système biométrique. Nous fournissons ensuite un score qui peut être utilisé pour paramétrer un système biométrique pour résister à cette attaque. Enfin, nous analysons les travaux qui présentent des scénarios de quasi-collision à l'aide de mesures telles que FMR, ceux qui fournissent la population critique pour un FMR donné, et ceux qui fournissent le FMR minimal pour des tailles de population données. Cette analyse nous permet de faire un retour critique et d'élargir les résultats de ces études.
- Section 2.4 : Cette section présente une analyse des fuites d'informations potentielles dans l'évaluation de la distance, particulièrement axée sur une distance offusquée à seuil. Les contributions comprennent une variété de scénarios de fuite d'informations, les attaques génériques correspondantes, leurs complexités et une correction d'un résultat présenté dans [83]. Les scénarios susmentionnés donnent lieu à de nouveaux scénarios d'attaques.
 1. Nous décrivons et étudions une nouvelle attaque, appelée attaque par accumulation, dans laquelle un serveur *honnête-mais-curieux* accumule des connaissances pendant l'authentification du client. Ce type d'attaque se produit lorsqu'il y a une fuite d'informations

mineure, mais non négligeable lors de l'évaluation de la distance.

2. Nous présentons de nouvelles stratégies d'attaque par des clients malveillants qui exploitent différents niveaux de fuites d'informations lors de l'évaluation de la distance. Nos résultats de complexité, qui détaillent le coût de ces attaques, s'appliquent à la fois aux attaques de type *recherche exhaustive hors ligne* qui exploitent une base de données divulguée (mais obscurcie) et aux attaques de type *recherche exhaustive en ligne* qui impliquent des interactions directes avec le serveur.
- Section 2.5 : Les attaques présentées sont basées sur une recherche exhaustive (*i.e.*, attaques par force brute) et ne nécessitent qu'une fuite minimale d'informations, à savoir un bit d'informations sur le succès de l'usurpation d'identité (*i.e.* 1 pour oui et 0 pour non). Ces attaques sont donc possibles quel que soit le schéma de transformation, le protocole ou la modalité biométrique utilisés. Nous supposons que le système biométrique utilise au mieux l'espace métrique sous-jacent afin de fournir les bornes théoriques les plus larges possibles sur la complexité des attaques par recherche exhaustive. Nous utilisons une modélisation probabiliste pour présenter deux scénarios d'attaque d'usurpation d'identité avec les bornes de sécurité associées et nous discutons de la sécurité d'une base de données de gabarits. Le premier, appelé "scénario de l'attaquant extérieur", illustre le cas où un individu non enregistré dans un service tente de se faire passer pour un utilisateur non spécifique de ce service. Plus précisément, nous envisageons la possibilité qu'un attaquant adapte sa stratégie de manière séquentielle. Le second scénario, appelé "scénario des attaquants extérieurs", englobe le cas où plusieurs attaquants attaquent le service en parallèle. Les bornes sur la complexité des attaques non ciblées fournissent la sécurité maximale réalisable. Enfin, nous formulons des recommandations concernant les paramètres de sécurité lors de la mise au point d'un système de reconnaissance.

Une grande partie de ces résultats provient de Durbet *et al.* [30, 18, 19].

Chapitre 3

Résumé : Dans ce chapitre, nous nous concentrons sur l'utilisation sécurisée des données biométriques. Plus précisément, nous visons à éviter ou limiter les attaques par recherche exhaustive hors ligne, car celles-ci sont plus dévastatrices que les attaques en ligne. En effet, les attaques en ligne sont généralement détectées ou ralenties par des contre-mesures telles que l'augmentation des délais entre les tentatives. Notre objectif est de construire un protocole cryptographique permettant d'authentifier un utilisateur uniquement à partir de ses données biométriques tout en se protégeant contre les attaques hors ligne. Nous procédons en deux étapes, d'abord en développant une nouvelle structure mathématique, puis en concevant deux protocoles distincts basés sur cette structure.

Premièrement, nous avons développé une nouvelle structure mathématique pour évaluer la distance entre deux vecteurs de manière homomorphique. Cette structure permet également de corriger les données fournies pour qu'elles correspondent aux données de référence si la distance est inférieure à un seuil donné. Les principaux avantages de cette construction résident dans le fait que la récupération des données de référence et lier deux enregistrements sont tous deux computationnellement difficiles.

Ensuite, avec cette nouvelle construction, nous avons développé deux protocoles ayant des objectifs distincts. Le premier protocole répond entièrement à notre problème initial en utilisant un mécanisme de preuve à *zéro connaissance* pour garantir qu'un client malveillant non limité par les capacités de calcul ne puisse obtenir d'informations utiles lors de l'authentification. De plus, nous proposons un deuxième protocole qui assure à la fois l'authentification et la reconstruction du secret original. Il convient de noter que ce protocole n'est approprié que si l'adversaire est supposé être borné polynomialement en termes de capacité de calcul, ce qui implique un niveau de sécurité inférieur. Cependant, les deux versions garantissent un haut niveau de sécurité contre les attaques hors ligne, car soit elles sont impossibles, soit elles ont un coût exponentiel.

Contexte : Un protocole d'authentification biométrique typique est composé de deux phases comprenant une série d'échanges entre le client et le serveur. Les phases d'inscription (*i.e.*, l'enrôlement) et de vérification (*i.e.*, le log-in) sont typiques d'un protocole d'authentification et donc *a fortiori*

d'un protocole d'authentification biométrique. Au cours de la phase d'inscription, des échantillons biométriques sont collectés auprès de l'utilisateur, puis traités pour obtenir un gabarit biométrique à l'aide d'un algorithme d'extraction de caractéristiques et d'une transformation. Le serveur stocke le gabarit (ou des informations dérivées du modèle, par exemple une clé cryptographique) avec l'identifiant de l'utilisateur. Dans la phase de vérification, l'utilisateur régénère son gabarit biométrique et l'utilise dans le protocole pour prouver au serveur son identité.

Le traitement, la transmission et le stockage d'informations dérivées des données biométriques tels que les gabarits constituent un aspect fondamental des applications biométriques en ligne. Les gabarits biométriques servent de données de référence principales pour reconnaître les individus de manière unique dans les applications et sont considérés comme faisant partie des informations personnelles. Par conséquent, la protection des informations biométriques des individus et de leur vie privée est d'une importance capitale dans les systèmes et applications biométriques. Certaines réglementations, telles que le règlement général sur la protection des données (RGPD) dans l'UE et le California Consumer Privacy Act (CCPA) pour les résidents de Californie aux États-Unis, ont été adoptées pour protéger les données et la vie privée des individus. La nature intrinsèquement variable des données biométriques, associée au fait que les caractéristiques biométriques d'un individu ne sont pas facilement renouvelables, rend la conception de protocoles d'authentification sécurisés basés sur la biométrie plus difficile et plus critique que la conception de protocoles basés sur des jetons (*i.e.*, des tokens) ou des mots de passe. Les efforts de recherche et de normalisation [23, 1, 2] ont permis d'identifier plusieurs exigences en matière de sécurisation des informations et des gabarits biométriques (voir la section 1.3.1 pour de plus amples détails).

Les sketches sécurisés sont des primitives cryptographiques utilisées pour protéger les gabarits biométriques avec des garanties de sécurité formelles. De manière informelle, étant donné un vecteur d'entrée biométrique x , un sketch S est dérivé par le biais d'un processus aléatoire. Le sketch S doit être *irréversible*, mais il doit aussi permettre de retrouver x en présence d'une autre donnée biométrique y proche de x vis-à-vis d'un seuil donnée. Dans ce chapitre, nous présentons une méthode générique permettant de construire une nouvelle famille de sketches sécurisés à partir d'une famille particulière de groupes. Nous appelons ces sketches sécurisés "Subset Product Sketch" (SPS) car la construction repose sur la multiplication d'éléments de groupe à partir d'un sous-ensemble d'éléments particulier que l'on appelle base. Les sketches sécurisés sont conçus pour permettre la correction de données variables dans le temps en cachant la donnée originelle que l'on voudrait récupérer. Cette capacité est très importante pour les applications où l'exactitude des données ne peut être garantie en raison de variations naturelles, comme dans les systèmes biométriques. Traditionnellement, les sketches sécurisés sont construits à l'aide de codes correcteurs d'erreurs pour traiter efficacement ces variations. En outre, les principes de la théorie de l'information garantissent la sécurité de ces sketches en gérant le compromis entre la *recupérabilité* (c'est-à-dire la capacité de correction) et la *confidentialité* des données. Plus précisément, les informations révélées par la publication de la valeur d'un sketch doivent être bornées vis-à-vis de la *min-entropie* [96]. Nous avons évité d'associer des notions de sécurité basées sur l'entropie aux schémas de sketches et fondons notre sécurité sur la difficulté des problèmes décisionnels et calculatoires.

Contributions détaillées : Dans ce chapitre, nous montrons comment construire un "Subset Product Sketch" (SPS) et comment l'utiliser pour mettre en place un secure sketch distant. SPS fonctionne comme un secure sketch, c'est-à-dire qu'il permet de retrouver un secret stocké en fournissant un secret similaire vis-à-vis d'une distance et d'un seuil, à l'aide de données d'aide. Ces données d'aide, générées à partir du secret, ne doivent révéler aucune information sur le secret lui-même. De plus, SPS doit satisfaire les propriétés d'être réutilisable et irréversible. Réutilisable signifie que si plusieurs données d'aide sont dérivées du même secret, les collecter ne donne pas plus d'informations sur le secret qu'une seule donnée. Irréversible signifie que, étant donné les données d'aide, il est difficile ou impossible de retrouver le secret original.

Pour ce faire, nous définissons d'abord une famille de groupes avec une propriété de factorisation unique. Cette propriété permet de cacher un secret dans l'exposant d'une base et, sous certaines conditions, de le retrouver de manière unique. Plus précisément, sous l'Heuristique Gaussienne, cette propriété permet de corriger de manière unique les erreurs sous le seuil de vérification. De plus, le fait de cacher le secret dans l'exposant assure l'irréversibilité sous des hypothèses computationnelles difficiles telles que le Computational Subset Product Problem (CSPP). Cette construction

permet de choisir de manière aléatoire un groupes et une base, assurant la réutilisabilité sous des problèmes décisionnels tels que le Decisional Subset Product Problem (DSPP). Ainsi, grâce à cette construction, SPS peut être implémenté et garantir toutes les propriétés mentionnées. Nous soulignons que notre SPS peut créer des sketches à partir de vecteurs d'entiers et tolérer une fraction linéaire d'erreurs par rapport à la distance L_1 (c'est-à-dire, pour $x, y \in \mathbb{Z}^n$, $L_1(x, y) = \sum_{i=0}^n |x_i - y_i|$). Cette fonctionnalité dépasse certaines constructions précédentes, qui sont limitées aux vecteurs binaires et à la distance de Hamming [63, 62, 60, 33].

Enfin, nous démontrons comment construire des secure sketches distants et des fuzzy extractors distants basés sur SPS. Notre premier protocole est un protocole de connaissance nulle, utilisant la vérifiabilité de la correction de la factorisation. Cette construction assure une résistance totale aux attaques hors ligne, car un client malveillant n'obtient aucune cible pour ses attaques. Dans le deuxième protocole, nous permettons au client de corriger sa donnée tout en assurant une résistance aux attaques hors ligne. Dans ce cas, comme le client reçoit des données d'aide obfusquées, la complexité de la recherche exhaustive est garantie par les hypothèses computationnelles mentionnées précédemment et la difficulté de casser l'obfuscation.

References

- [1] ISO/IEC24745 :2011 : *Information technology – Security techniques – Biometric information protection*. Standard. International Organization for Standardization, 2011.
- [2] ISO/IEC30136 :2018(E) : *Information technology – Performance testing of biometric template protection scheme*. Standard. International Organization for Standardization, 2018.
- [3] APPLE. *About Face ID advanced technology*. <https://support.apple.com/en-us/102381>.
- [4] APPLE. *About Touch ID advanced security technology*. <https://support.apple.com/en-us/HT204587>.
- [5] Chris BURT. *Palm vein biometrics deployed for payments at Paris grocery store*. <https://www.biometricupdate.com/202407/palm-vein-biometrics-deployed-for-payments-at-paris-grocery-store>.
- [6] FVC-ONGOING. *Published Results*. <https://biolab.csr.unibo.it/fvcongoing/UI/Form/PublishedAlgs.aspx>.
- [7] GOOGLE. *Measuring Biometric Unlock Security*. <https://source.android.com/docs/security/features/biometric/measure>.
- [8] Didier GUILLERM. *History of biometrics*. <https://www.biometrie-online.net/biometrie/histoire>.
- [9] Commission Nationale de l'Informatique et des LIBERTÉS. *Biométrie*. <https://www.cnil.fr/fr/biometrie>.
- [10] Ahmed SHERIF. *Statistica : Global market share held by operating systems for desktop PCs, from January 2013 to February 2024*. <https://www.statista.com/statistics/218089/global-market-share-of-windows-7/>.
- [11] John DAUGMAN. « Understanding Biometric Entropy and Iris Capacity : Avoiding Identity Collisions on National Scales ». In : *Adv. Artif. Intell. Mach. Learn.* 4.2 (2024), p. 2152-2163. DOI : 10.54364/AAIML.2024.42123. URL : <https://doi.org/10.54364/aaiml.2024.42123>.
- [12] Axel DURBET, Koray KARABINA and Kevin THIRY-ATIGHEHCHI. *Generic Construction of Secure Sketches from Groups*. Cryptology ePrint Archive, Paper 2024/1224. <https://eprint.iacr.org/2024/1224>. 2024. URL : <https://eprint.iacr.org/2024/1224>.
- [13] Torbjörn GRANLUND, Gunnar SJÖDIN, Hans RIESEL, Richard STALLMAN, Brian BEUNING, Doug LEA, John AMANATIDES, Paul ZIMMERMANN, Ken WEBER, Per BOTHNER, Joachim HOLLMAN, Bennet YEE, Andreas SCHWAB, Robert HARLEY, David SEAL, Torsten EKEDAHL, Linus NORDBERG, Kevin RYDE, Kent BOORTZ, Steve ROOT, Gerardo BALLABIO, Jason MOXHAM, Pedro GIMENO, Niels MÖLLER, Alberto ZANONI, Marco BODRATO, Marco BODRATO, David HARVEY, Martin BOIJ, Marc GLISSE, David S MILLER, Mark SOFRONIOU and Ulrich WEIGAND. *The GNU Multiple Precision Arithmetic Library*. <https://gmp.lib.org/>. 2024. URL : <https://gmp.lib.org/>.
- [14] Patrick GROTHOR, Mei NGAN, Kayee HANAOKA, Joyce C. YANG and Austin HOM. *Face Recognition Technology Evaluation (FRTE), Part 1 : Verification*. NIST Interagency Report. 2024.
- [15] Damodaran HARIKRISHNAN, Sunil KUMAR, Shelbi JOSEPH and Kishor Krishnan NAIR. « Towards a fast and secure fingerprint authentication system based on a novel encoding scheme ». In : *International Journal of Electrical Engineering & Education* 61.1 (2024), p. 100-112.

- [16] Victor SHOUP. *NTL : A Library for doing Number Theory*. <https://libnt1.org/>. 2024. URL : <https://libnt1.org/>.
- [17] Daniel APON, Chloe CACHET, Benjamin FULLER, Peter HALL and Feng-Hao LIU. « Nonmal-leable digital lockers and robust fuzzy extractors in the plain model ». In : *Advances in Cryptology-ASIACRYPT 2022 : 28th International Conference on the Theory and Application of Cryptology and Information Security, Taipei, Taiwan, December 5-9, 2022, Proceedings, Part IV*. Springer. 2023, p. 353-383.
- [18] Axel DURBET, Paul-Marie GROLLEMUND and Kevin THIRY-ATIGHEHCHI. *Untargeted Near-collision Attacks in Biometric Recognition*. 2023. arXiv : 2304.01580 [cs.CR].
- [19] Axel DURBET, Kevin THIRY-ATIGHEHCHI, Dorine CHAGNON and Paul-Marie GROLLEMUND. *Exploit the Leak : Understanding Risks in Biometric Matchers*. 2023. arXiv : 2307.13717 [cs.CR].
- [20] Venkatesan GURUSWAMI, Atri RUDRA and Madhu SUDAN. *Essential Coding Theory*. 2023. URL : <https://cse.buffalo.edu/faculty/atricourses/coding-theory/book/web-coding-book.pdf>.
- [21] Kaiji MOTEKI and Soma HAYASHI. « A groupwise approach to the birthday paradox ». In : *Available at SSRN 4593602* (2023).
- [22] Somnath PANJA, Nikita TRIPATHI, Shaoquan JIANG and Reihaneh SAFAVI-NAINI. *Robust and Reusable Fuzzy Extractors and their Application to Authentication from Iris Data*. Cryptology ePrint Archive, Paper 2023/284. 2023.
- [23] Shreyansh SHARMA, Anil SAINI and Santanu CHAUDHURY. « A survey on biometric cryptosystems and their applications ». In : *Computers & Security* (2023), p. 103458.
- [24] THALES. *The History of Biometric Authentication*. <https://www.thalesgroup.com/en/markets/digital-identity-and-security/government/inspired/history-of-biometric-authentication>. 2023.
- [25] Esteban VÁZQUEZ and Artur COSTA-PAZO. *Defining the core accuracy metrics of biometric systems*. <https://alicebiometrics.com/en/defining-the-core-accuracy-metrics-of-biometric-systems/>. Security agency blog. 2023.
- [26] Zaynab ALMUTAIRI and Hebah ELGIBREEN. « A review of modern audio deepfake detection methods : Challenges and future directions ». In : *Algorithms* 15.5 (2022), p. 155.
- [27] Emilio ANDREOZZI, Riccardo SABBADINI, Jessica CENTRACCHIO, Paolo BIFULCO, Andrea IRACE, Giovanni BREGGIO and Michele RICCIO. « Multimodal Finger Pulse Wave Sensing : Comparison of Forcecardiography and Photoplethysmography Sensors ». In : *Sensors* 22.19 (2022), p. 7566.
- [28] Sílvia CASACUBERTA, Julia HESSE and Anja LEHMANN. « SoK : Oblivious pseudorandom functions ». In : *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*. IEEE. 2022, p. 625-646.
- [29] Axel DURBET, Paul-Marie GROLLEMUND, Pascal LAFOURCADE, Denis MIGDAL and Kevin THIRY-ATIGHEHCHI. « Authentication Attacks on Projection-based Cancelable Biometric Schemes ». In : *Proceedings of the 19th International Conference on Security and Cryptography, SECRYPT 2022, Lisbon, Portugal, July 11-13, 2022*. Sous la dir. de Sabrina De Capitani di VIMERCATI and Pierangela SAMARATI. SCITEPRESS, 2022, p. 568-573.
- [30] Axel DURBET, Paul-Marie GROLLEMUND, Pascal LAFOURCADE and Kevin THIRY-ATIGHEHCHI. « Near-collisions and Their Impact on Biometric Security ». In : *Proceedings of the 19th International Conference on Security and Cryptography, SECRYPT 2022, Lisbon, Portugal, July 11-13, 2022*. Sous la dir. de Sabrina De Capitani di VIMERCATI and Pierangela SAMARATI. SCITEPRESS, 2022, p. 382-389.
- [31] Tanguy GERNOT and Patrick LACHARME. « Biometric masterkeys ». In : *Computers & Security* 116 (2022), p. 102642. ISSN : 0167-4048. DOI : 10.1016/j.cose.2022.102642. URL : <https://www.sciencedirect.com/science/article/pii/S0167404822000414>.

- [32] Arun ROSS, Sudipta BANERJEE and Anurag CHOWDHURY. « Deducing health cues from biometric data ». In : *Computer Vision and Image Understanding* 221 (2022), p. 103438. ISSN : 1077-3142.
- [33] Ran CANETTI, Benjamin FULLER, Omer PANETH, Leonid REYZIN and Adam SMITH. « Reusable fuzzy extractors for low-entropy distributions ». In : *Journal of Cryptology* 34 (2021), p. 1-33.
- [34] William FISHER, Don FAATZ, Mark RUSSELL, Christopher BROWN, Sanjeev SHARMA, Sudhi UMARJI and Karen SCARFONE. *Using Mobile Device Biometrics for Authenticating First Responders*. National Institute of Standards and Technology. 2021.
- [35] Rohit MUTALIK, Dhairya CHHEDA, Zeeshan SHAIKH and Dhanashree TORADMALLE. *Rockyou*. 2021. DOI : 10.21227/gzcg-yc14. URL : <https://dx.doi.org/10.21227/gzcg-yc14>.
- [36] Pascal Sungu NGOY, Kaninda MUSUMBU and Duncan Kioi GATHUNGU. « Combining epidemic model and deep learning to study cyber attacks ». In : *2021 International Conference on Computer Communication and Artificial Intelligence (CCAI)*. IEEE. 2021, p. 150-154.
- [37] Cornelia OTT, Sven PUCHINGER and Martin BOSSERT. « Bounds and genericity of sum-rank-metric codes ». In : *2021 XVII International Symposium "Problems of Redundancy in Information and Control Systems"(REDUNDANCY)*. IEEE. 2021, p. 119-124.
- [38] Shoukat ALI, Koray KARABINA and Emrah KARAGOZ. « Formal Accuracy Analysis of a Biometric Data Transformation and Its Application to Secure Template Generation ». In : *SECRYPT 2020*. ScitePress, 2020, p. 485-496.
- [39] Gabrielle DE MICHELI, Pierrick GAUDRY and Cécile PIERROT. « Asymptotic complexities of discrete logarithm algorithms in pairing-relevant finite fields ». In : *Advances in Cryptology—CRYPTO 2020 : 40th Annual International Cryptology Conference, CRYPTO 2020, Santa Barbara, CA, USA, August 17–21, 2020, Proceedings, Part II* 40. Springer. 2020, p. 32-61.
- [40] Benjamin FULLER, Xianrui MENG and Leonid REYZIN. « Computational fuzzy extractors ». In : *Information and Computation* 275 (2020), p. 1-20.
- [41] Alexander HELM and Alexander MAY. « The power of few qubits and collisions—subset sum below Grover's bound ». In : *Post-Quantum Cryptography : 11th International Conference, PQCrypto 2020, Paris, France, April 15–17, 2020, Proceedings*. Springer. 2020, p. 445-460.
- [42] MANISHA and Nitin KUMAR. « Cancelable biometrics : a comprehensive survey ». In : *Artificial Intelligence Review* 53.5 (2020), p. 3403-3446.
- [43] Junichi TOMIDA. « Tightly secure inner product functional encryption : Multi-input and function-hiding constructions ». In : *Theoretical Computer Science* 833 (2020), p. 56-86.
- [44] Junichi TOMIDA and Katsuyuki TAKASHIMA. « Unbounded inner product functional encryption from bilinear maps ». In : *Japan Journal of Industrial and Applied Mathematics* 37.3 (2020), p. 723-779.
- [45] Asma AL KINDI, Dawood AL ABRI, Ahmed AL MAASHRI and Fahad BAIT-SHIGINAH. « Analysis of malware propagation behavior in Social Internet of Things ». In : *International Journal of Communication Systems* 32.15 (2019), e4102.
- [46] Yoshinori AONO, Thomas ESPITAU and Phong Q NGUYEN. « Random lattices : theory and practice ». In : *Preprint*. https://espitau.github.io/bin/random_lattice.pdf (2019).
- [47] Kathleen BRUSH, Nabil EL ACHRAOUI, Jennifer BOYD, Jacob JOHNSON, Randy CHEPENIK, Tarik MCLEAN, Sadida SIDDIQUI, Aditee VERMA, John SHERIDAN, Avery LEIDER and C. Charles TAPPERT. « Index of Difficulty Measurement for Handedness with Biometric Authentication ». In : *HCI International 2019 – Late Breaking Posters*. Sous la dir. de Constantine STEPHANIDIS and Margherita ANTONA. Cham : Springer International Publishing, 2019, p. 413-423.
- [48] Naser DAMER, Fadi BOUTROS, Alexandra Moseguí SALADIÉ, Florian KIRCHBUCHNER and Arjan KUIJPER. « Realistic Dreams : Cascaded Enhancement of GAN-generated Images with an Example in Face Morphing Attacks ». In : *2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. 2019, p. 1-10.

- [49] D. Steven GALBRAITH and Lukas ZOBERNIG. « Obfuscated Fuzzy Hamming Distance and Conjunctions from Subset Product Problems ». In : *Theory of Cryptography*. Sous la dir. de Dennis HOFHEINZ and Alon ROSEN. Cham : Springer International Publishing, 2019, p. 81-110.
- [50] Anas HUSSEIS, Judith LIU-JIMENEZ, Ines GOICOECHEA-TELLERIA and Raul SANCHEZ-REILLO. « A Survey in Presentation Attack and Presentation Attack Detection ». In : *2019 International Carnahan Conference on Security Technology (ICCST)*. 2019, p. 1-13.
- [51] Ulrich SCHERHAG, Christian RATHGEB, Johannes MERKLE, Ralph BREITHAUPT and Christoph BUSCH. « Face Recognition Systems Under Morphing Attacks : A Survey ». In : *IEEE Access* 7 (2019), p. 23012-23026.
- [52] Mika WESTERLUND. « The emergence of deepfake technology : A review ». In : *Technology innovation management review* 9.11 (2019).
- [53] Wencheng YANG, Song WANG, Jiankun HU, Guanglou ZHENG and Craig VALLI. « Security and Accuracy of Fingerprint-Based Biometrics : A Review ». In : *Symmetry* 11.2 (2019). ISSN : 2073-8994.
- [54] Jung Hee CHEON, Jinhyuck JEONG, Dongwoo KIM and Jongchan LEE. « A Reusable Fuzzy Extractor with Practical Storage Size : Modifying Canetti et al.'s Construction ». In : *Information Security and Privacy*. 2018, p. 28-44.
- [55] David EVANS, Vladimir KOLESNIKOV, Mike ROSULEK et al. « A pragmatic introduction to secure multi-party computation ». In : *Foundations and Trends® in Privacy and Security* 2.2-3 (2018), p. 70-246.
- [56] Miguel A FERRER, Moises DIAZ, Cristina CARMONA-DUARTE and Réjean PLAMONDON. « A biometric attack case based on signature synthesis ». In : *2018 International Carnahan Conference on Security Technology (ICCST)*. IEEE. 2018, p. 1-6.
- [57] Olegs NIKISINS, Teodors EGLITIS, André ANJOS and Sébastien MARCEL. « Fast cross-correlation based wrist vein recognition algorithm with rotation and translation compensation ». In : *2018 International Workshop on Biometrics and Forensics (IWBF)*. IEEE. 2018, p. 1-7.
- [58] W. George QUINN, Patrick GROTHOR and James MATEY. *IREX IX part one : Performance of Iris recognition algorithms*. US Department of Commerce, National Institute of Standards and Technology ..., 2018.
- [59] David G. REES. *Essential statistics*. Chapman and Hall/CRC, 2018.
- [60] Yunhua WEN and Shengli LIU. « Reusable Fuzzy Extractor from LWE ». In : *Information Security and Privacy*. 2018, p. 13-27.
- [61] Yunhua WEN and Shengli LIU. « Robustly Reusable Fuzzy Extractor from Standard Assumptions ». In : *Advances in Cryptology – ASIACRYPT 2018*. 2018, p. 459-489.
- [62] Yunhua WEN, Shengli LIU and Shuai HAN. « Reusable fuzzy extractor from the decisional Diffie–Hellman assumption ». In : t. 86. 2018, p. 2495-2512.
- [63] Daniel APON, Chongwon CHO, Karim ELDEFRAWY and Jonathan KATZ. « Efficient, reusable fuzzy extractors from LWE ». In : *Cyber Security Cryptography and Machine Learning : First International Conference, CSCML 2017, Beer-Sheva, Israel, June 29-30, 2017, Proceedings 1*. Springer. 2017, p. 1-18.
- [64] Yerra Shankar RAO, Aswin Kumar RAUTA, Hemraj SAINI and Tarini Charana PANDA. « Mathematical model for cyber attack in computer network ». In : *International Journal of Business Data Communications and Networking (IJBDN)* 13.1 (2017), p. 58-65.
- [65] Sheldon M. Ross. *Introductory statistics*. Academic Press, 2017.
- [66] Mulagala SANDHYA and V.N.K. Munaga PRASAD. « Biometric template protection : A systematic literature review of approaches and modalities ». In : *Biometric Security and Privacy : Opportunities & Challenges in The Big Data Era* (2017), p. 323-370.
- [67] Ran CANETTI, Benjamin FULLER, Omer PANETH, Leonid REYZIN and Adam D. SMITH. « Reusable Fuzzy Extractors for Low-Entropy Distributions ». In : *Advances in Cryptology - EUROCRYPT 2016*. 2016, p. 117-146.

- [68] Hao CHEN. *A Measure Version of Gaussian Heuristic*. Cryptology ePrint Archive, Paper 2016/439. 2016. URL : <https://eprint.iacr.org/2016/439>.
- [69] Ivana CHINGOVSKA, Nesli ERDOGMUS, André ANJOS and Sébastien MARCEL. « Face recognition systems under spoofing attacks ». In : *Face Recognition Across the Imaging Spectrum*. Springer, 2016, p. 165-194.
- [70] John DAUGMAN and Cathryn DOWNING. « Searching for doppelgängers : Assessing the universality of the IrisCode impostors distribution ». In : *IET Biometrics* 5.2 (2016), p. 65-75.
- [71] Koray KARABINA and Onur CANPOLAT. « A new cryptographic primitive for noise tolerant template security ». In : *Pattern Recognition Letters* 80 (2016), p. 70-75. ISSN : 0167-8655.
- [72] Pecatte TIMOTHÉE and Somindu Chaya RAMANNA. *Tutorial 10 for Information Theory*. 2016. URL : https://perso.ens-lyon.fr/omar.fawzi/teaching/it/tutorial10_questions.pdf.
- [73] Adam CZAJKA. « Pupil dynamics for iris liveness detection ». In : *IEEE Transactions on Information Forensics and Security* 10.4 (2015), p. 726-735.
- [74] John DAUGMAN. « Information theory and the iriscode ». In : *IEEE transactions on information forensics and security* 11.2 (2015), p. 400-409.
- [75] J. Morris DWORKIN. « SHA-3 standard : Permutation-based hash and extendable-output functions ». In : (2015).
- [76] V. Oleg KOMOGORTSEV, Alexey KARPOV and D. Corey HOLLAND. « Attack of Mechanical Replicas : Liveness Detection With Eye Movements ». In : *IEEE Transactions on Information Forensics and Security* 10.4 (2015), p. 716-725.
- [77] M. Vishal PATEL, K. Nalini RATHA and Rama CHELLAPPA. « Cancelable Biometrics : A review ». In : *IEEE Signal Processing Magazine* 32.5 (2015), p. 54-65.
- [78] Pedro TOME and Sébastien MARCEL. « On the vulnerability of palm vein recognition to spoofing attacks ». In : *2015 International Conference on Biometrics (ICB)*. 2015, p. 319-325.
- [79] John AHLGREN. « The probability distribution for draws until first success without replacement ». In : (2014).
- [80] Julien BRINGER, Herve CHABANNE, Melanie FAVRE, Alain PATEY, Thomas SCHNEIDER and Michael ZOHNER. « GSHADE : faster privacy-preserving distance computation and biometric identification ». In : *Proceedings of the 2nd ACM workshop on Information hiding and multimedia security*. 2014, p. 187-198.
- [81] Marco FERRANTE and Monica SALTALAMACCHIA. « The Coupon Collector's problem ». In : *MATerials MATemàtics* (mai 2014), p. 35.
- [82] Matteo FERRARA, Annalisa FRANCO and Davide MALTONI. « The magic passport ». In : *IJCB 2014 - 2014 IEEE/IAPR International Joint Conference on Biometrics* (déc. 2014).
- [83] Elena PAGNIN, Christos DIMITRAKAKIS, Aysajan ABIDIN and Aikaterini MITROKOTSA. « On the Leakage of Information in Biometric Authentication ». In : *Progress in Cryptology - INDOCRYPT 2014 - 15th International Conference on Cryptology in India, New Delhi, India, December 14-17, 2014, Proceedings*. Sous la dir. de Willi MEIER and Debdeep MUKHOPADHYAY. T. 8885. Lecture Notes in Computer Science. Springer, 2014, p. 265-280.
- [84] Rima BELGUECHI, Estelle CHERRIER, Christophe ROSENBERGER and Samy AIT-AOUDIA. « Operational bio-hash to preserve privacy of fingerprint minutiae templates ». In : *IET biometrics* 2.2 (2013), p. 76-84.
- [85] J. Daniel BERNSTEIN, Stacey JEFFERY, Tanja LANGE and Alexander MEURER. « Quantum algorithms for the subset-sum problem ». In : *Post-Quantum Cryptography : 5th International Workshop, PQCrypto 2013, Limoges, France, June 4-7, 2013. Proceedings* 5. Springer. 2013, p. 16-33.
- [86] Bin DAI, Shilin DING and Grace WAHBA. « Multivariate bernoulli distribution ». In : (2013).
- [87] Melissa GYMREK, Amy L. MCGUIRE, David GOLAN, Eran HALPERIN and Yaniv ERLICH. « Identifying Personal Genomes by Surname Inference ». In : *Science* 339.6117 (2013), p. 321-324.

- [88] MIT News Office LARRY HARDESTY. *How hard is it to 'de-anonymize' cellphone data?* <https://news.mit.edu/2013/how-hard-it-de-anonymize-cellphone-data>. 2013.
- [89] Sancheng PENG, Shui YU and Aimin YANG. « Smartphone malware and its propagation modeling : A survey ». In : *IEEE Communications Surveys & Tutorials* 16.2 (2013), p. 925-941.
- [90] Reihaneh SAFAVI-NAINI and Ran CANETTI. *Advances in Cryptology—CRYPTO 2012 : 32nd Annual Cryptology Conference, Santa Barbara, CA, USA, August 19-23, 2012, Proceedings*. T. 7417. Springer, 2012.
- [91] Koen SIMOENS, Julien BRINGER, Hervé CHABANNE and Stefaan SEYS. « A framework for analyzing template security and privacy in biometric authentication systems ». In : *IEEE Transactions on Information forensics and security* 7.2 (2012), p. 833-841.
- [92] Yuanmi CHEN and Phong Q NGUYEN. « BKZ 2.0 : Better lattice security estimates ». In : *International Conference on the Theory and Application of Cryptology and Information Security*. Springer. 2011, p. 1-20.
- [93] Won KIM, Ok-Ran JEONG, Chulyun KIM and Jungmin So. « The dark side of the Internet : Attacks, costs and responses ». In : *Information systems* 36.3 (2011), p. 675-705.
- [94] H. Bernhard KORTE and Jens VYGEN. *Combinatorial optimization*. T. 1. Springer, 2011.
- [95] J. Pedro PEDROSO. « Optimization with Gurobi and Python ». In : *INESC Porto and Universidade do Porto, Porto, Portugal 1* (2011).
- [96] Leonid REYZIN. « Some Notions of Entropy for Cryptography : (Invited Talk) ». In : *International Conference on Information Theoretic Security*. Springer. 2011, p. 138-142.
- [97] Zachariah DIETZ and Michael E. SCHUCKERS. « A central-limit theorem for a single-false match rate ». In : *Biometric Technology for Human Identification VII*. Sous la dir. de B. V. K. Vijaya KUMAR, Salil PRABHAKAR and Arun A. ROSS. T. 7667. International Society for Optics and Photonics. SPIE, 2010, 76670F. DOI : 10.1117/12.849746. URL : <https://doi.org/10.1117/12.849746>.
- [98] W Cary HUFFMAN and Vera PLESS. *Fundamentals of error-correcting codes*. Cambridge university press, 2010.
- [99] Bin MA and Xiaoming SUN. « More efficient algorithms for closest string and substring problems ». In : *SIAM Journal on Computing* 39.4 (2010), p. 1432-1443.
- [100] Osama OUDA, Norimichi TSUMURA and Toshiya NAKAGUCHI. « Bioencoding : A reliable tokenless cancelable biometrics scheme for protecting iriscodes ». In : *IEICE TRANSACTIONS on Information and Systems* 93.7 (2010), p. 1878-1888.
- [101] Feng QI. « Bounds for the ratio of two gamma functions ». In : *Journal of Inequalities and Applications* 2010 (2010), p. 1-84.
- [102] Michael E SCHUCKERS. *Computational methods in biometric authentication : statistical methods for performance evaluation*. Springer Science & Business Media, 2010.
- [103] Petra BERENBRINK and Thomas SAUERWALD. « The Weighted Coupon Collector's Problem and Applications ». In : *Computing and Combinatorics*. Sous la dir. de Q. Hung Ngo. Berlin, Heidelberg : Springer Berlin Heidelberg, 2009, p. 449-458.
- [104] John DAUGMAN. « How iris recognition works ». In : *The essential guide to image processing*. Elsevier, 2009, p. 715-739.
- [105] Ayman JARROUS and Benny PINKAS. « Secure hamming distance based computation and its applications ». In : *Applied Cryptography and Network Security : 7th International Conference, ACNS 2009, Paris-Rocquencourt, France, June 2-5, 2009. Proceedings 7*. Springer. 2009, p. 107-124.
- [106] Robert KONIG, Renato RENNER and Christian SCHAFFNER. « The Operational Meaning of Min- and Max-Entropy ». In : *IEEE Transactions on Information Theory* 55.9 (2009), p. 4337-4347. DOI : 10.1109/TIT.2009.2025545.
- [107] Oded REGEV. « On lattices, learning with errors, random linear codes, and cryptography ». In : *Journal of the ACM (JACM)* 56.6 (2009), p. 1-40.

- [108] Koen SIMOENS, Pim TUYLS and Bart PRENEEL. « Privacy Weaknesses in Biometric Sketches ». In : *2009 30th IEEE Symposium on Security and Privacy*. 2009, p. 188-203.
- [109] Yevgeniy DODIS, Rafail OSTROVSKY, Leonid REYZIN and Adam SMITH. « Fuzzy extractors : How to generate strong keys from biometrics and other noisy data ». In : *SIAM journal on computing* 38.1 (2008), p. 97-139.
- [110] Nicolas GAMA and Q. Phong NGUYEN. « Predicting Lattice Reduction ». In : *Advances in Cryptology – EUROCRYPT 2008*. 2008, p. 31-51.
- [111] Joel RATSABY. « Estimate of the number of restricted integer-partitions ». In : *Applicable Analysis and Discrete Mathematics* (2008), p. 222-233.
- [112] Virginia RUIZ-ALBACETE, Pedro TOME-GONZALEZ, Fernando ALONSO-FERNANDEZ, Javier GALBALLY, Julian FIERREZ and Javier ORTEGA-GARCIA. « Direct attacks using fake images in iris verification ». In : *European workshop on biometrics and identity management*. Springer. 2008, p. 181-190.
- [113] Qiang TANG, Julien BRINGER, Herve CHABANNE and David POINTCHEVAL. « A Formal Study of the Privacy Concerns in Biometric-Based Remote Authentication Schemes ». In : t. 4991. Avr. 2008.
- [114] Julian FIERREZ, Javier ORTEGA-GARCIA, Doroteo Torre TOLEDANO and Joaquin GONZALEZ-RODRIGUEZ. « BioSec baseline corpus : A multimodal biometric database ». In : *Pattern Recognition* 40.4 (2007), p. 1389-1392.
- [115] Arvind NARAYANAN and Vitaly SHMATIKOV. *How To Break Anonymity of the Netflix Prize Dataset*. 2007.
- [116] Ari JUELS. « A Fuzzy Vault scheme ». In : *Designs, Codes and Cryptography* 38 (fév. 2006), p. 237-257.
- [117] Phong Q NGUYEN and Damien STEHLÉ. « LLL on the Average ». In : *Proceedings of the 7th International Conference on Algorithmic Number Theory*. 2006, p. 238-256.
- [118] Andrzej PACUT and Adam CZAJKA. « Aliveness detection for iris biometrics ». In : *Proceedings 40th annual 2006 international carnahan conference on security technology*. IEEE. 2006, p. 122-129.
- [119] M. Cover THOMAS and A. Thomas JOY. *Elements of information theory*. Wiley-Interscience, 2006.
- [120] Xavier BOYEN, Yevgeniy DODIS, Jonathan KATZ, Rafail OSTROVSKY and Adam SMITH. « Secure remote authentication using biometric data ». In : *Advances in Cryptology–EUROCRYPT 2005 : 24th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Aarhus, Denmark, May 22-26, 2005. Proceedings 24*. Springer. 2005, p. 147-163.
- [121] Michael O RABIN. « How to exchange secrets with oblivious transfer ». In : *Cryptology ePrint Archive* (2005).
- [122] Xavier BOYEN. « Reusable cryptographic fuzzy extractors ». In : *Proceedings of the 11th ACM conference on Computer and Communications Security*. 2004, p. 82-91.
- [123] Yevgeniy DODIS, Leonid REYZIN and Adam D. SMITH. « Fuzzy Extractors : How to Generate Strong Keys from Biometrics and Other Noisy Data ». In : *Advances in Cryptology - EUROCRYPT 2004, International Conference on the Theory and Applications of Cryptographic Techniques, Interlaken, Switzerland, May 2-6, 2004, Proceedings*. 2004, p. 523-540.
- [124] Marcos FAUNDEZ-ZANUY. « On the vulnerability of biometric security systems ». In : *IEEE Aerospace and Electronic Systems Magazine* 19.6 (2004), p. 3-8.
- [125] Claudio MENESES, Zhaosong LU, Carlos OLIVEIRA and Panos PARDALOS. « Optimal Solutions for the Closest String Problem via Integer Programming ». In : *INFORMS Journal on Computing - INFORMS* 16 (nov. 2004), p. 419-429.
- [126] Victor SHOUP. *Sequences of games : a tool for taming complexity in security proofs*. Cryptology ePrint Archive, Paper 2004/332. <https://eprint.iacr.org/2004/332>. 2004. URL : <https://eprint.iacr.org/2004/332>.

- [127] John DAUGMAN. « The importance of being random : statistical principles of iris recognition ». In : *Pattern recognition* 36.2 (2003), p. 279-291.
- [128] A. Patricia EVANS and D. Andrew SMITH. « Complexity of approximating closest substring problems ». In : *Fundamentals of Computation Theory : 14th International Symposium, FCT 2003, Malmö, Sweden, August 12-15, 2003. Proceedings 14*. Springer. 2003, p. 210-221.
- [129] Ming LI, Bin MA and Lusheng WANG. « On the closest string and substring problems ». In : *Journal of the ACM (JACM)* 49.2 (2002), p. 157-171.
- [130] Tsutomu MATSUMOTO. « Gummy and conductive silicone rubber fingers importance of vulnerability analysis ». In : *International conference on the theory and application of cryptology and information security*. Springer. 2002, p. 574-575.
- [131] Tsutomu MATSUMOTO, Hiroyuki MATSUMOTO, Koji YAMADA and Satoshi HOSHINO. « Impact of artificial "gummy" fingers on fingerprint systems ». In : *Optical Security and Counterfeit Deterrence Techniques IV*. T. 4677. SPIE. 2002, p. 275-289.
- [132] Jens GRAMM, Rolf NIEDERMEIER and Peter ROSSMANITH. « Exact solutions for closest string and related problems ». In : *International Symposium on Algorithms and Computation*. Springer. 2001, p. 441-453.
- [133] K. Nalini RATHA, H. Jonathan CONNELL and M. Ruud BOLLE. « An analysis of minutiae matching strength ». In : *International Conference on Audio-and Video-Based Biometric Person Authentication*. Springer. 2001, p. 223-228.
- [134] Anil K. JAIN, Salil PRABHAKAR, L. HONG and Sharath PANKANTI. « Filterbank-based fingerprint matching ». In : *IEEE Transactions on Image Processing* 9.5 (2000), p. 846-859.
- [135] Ari JUELS and Martin WATTENBERG. « A Fuzzy Commitment Scheme ». In : *Proceedings of the 6th ACM Conference on Computer and Communications Security. CCS '99*. Kent Ridge Digital Labs, Singapore : Association for Computing Machinery, 1999, p. 28-36.
- [136] Ari JUELS and Martin WATTENBERG. « A Fuzzy Commitment Scheme ». In : *Proceedings of the 6th ACM Conference on Computer and Communications Security*. Association for Computing Machinery, 1999, p. 28-36.
- [137] Ari JUELS and Martin WATTENBERG. « A fuzzy commitment scheme ». In : *Proceedings of the 6th ACM conference on Computer and communications security*. 1999, p. 28-36.
- [138] Oded GOLDREICH. « Secure multi-party computation ». In : *Manuscript. Preliminary version* 78.110 (1998), p. 1-108.
- [139] Moti FRANCES and Amélie LITMAN. « On covering problems of codes ». In : *Theory of Computing Systems* 30.2 (avr. 1997), p. 113-119.
- [140] Bangalore S. MANJUNATH and Wei-Ying MA. « Texture features for browsing and retrieval of image data ». In : *IEEE Transactions on pattern analysis and machine intelligence* 18.8 (1996), p. 837-842.
- [141] Joe KILIAN. « Founding cryptography on oblivious transfer ». In : *Proceedings of the twentieth annual ACM symposium on Theory of computing*. 1988, p. 20-31.
- [142] Scott KIRKPATRICK, Daniel C. GELATT and Mario P. VECCHI. « Optimization by Simulated Annealing ». In : *Science* 220.4598 (1983), p. 671-680.
- [143] Vasek CHVATAL. « A greedy heuristic for the set-covering problem ». In : *Mathematics of operations research* 4.3 (1979), p. 233-235.
- [144] Daniel DEFAYS. « An efficient algorithm for a complete link method ». In : *The Computer Journal* 20.4 (jan. 1977), p. 364-366. ISSN : 0010-4620.
- [145] FJ MACWILLIAMS. « The theory of error-correcting codes ». In : *Elsevier Science Publishers BV* 2 (1977), p. 39-47.
- [146] WW PETERSON. « Error-correcting codes ». In : *Cambridge, MA : MIT Press* 2 (1972), p. 208-213.
- [147] Richard W HAMMING. « Error detecting and error correcting codes ». In : *The Bell system technical journal* 29.2 (1950), p. 147-160.

List of Figures

| | | |
|------|---|-----|
| 1.1 | Schematic view of a biometric system for enrollment and authentication. | 7 |
| 1.2 | Attack points in a generic biometric recognition system. | 9 |
| 1.3 | Classification of attacks on biometric cryptosystems (Sharma <i>et al.</i> [23]). | 12 |
| 2.1 | Representation of a database with 5 templates marked by a dot. | 21 |
| 2.2 | Representation of a ϵ -master-template marked by a cross for a database with 3 templates marked by a dot. | 22 |
| 2.3 | Representation of a clustered database for a threshold of ϵ | 23 |
| 2.4 | The figure depicts the representation of a clustered database for a threshold of ϵ . The ϵ -master template of each cluster is represented by a cross. | 24 |
| 2.5 | The figure depicts the representation of the database partitioning algorithm. | 29 |
| 2.6 | Illustration of near-collisions of two templates. | 32 |
| 2.7 | Difference between targeted and untargeted attacks from the attacker point of view. The highlighted surface represents the attackable area. | 43 |
| 2.8 | Exploiting the error position leaked in the case \mathbb{Z}_4^5 and the hidden vector or missing coordinates is $(0, 1, 3, 2, 2)$ | 50 |
| 2.9 | Representation of a κ -adaptative seeking a template of a user. The cross symbols indicate unsuccessful attempts, while the hashed areas represent the knowledge gained about the absence of templates in those regions. | 55 |
| 3.1 | Relations between the hardness of problems DSPP, CSPP, DLP, and the security properties reusability and irreversibility of SPS. | 76 |
| 3.2 | SPS-based Fuzzy Authentication Scheme | 91 |
| 3.3 | SPS-based Authentication. | 92 |
| 3.4 | Schematic view of SPS. | 92 |
| 3.5 | OSS scheme based on SPS | 94 |
| 3.6 | SPS-based Sketch. | 94 |
| 3.7 | SPS-based <i>rRec</i> | 95 |
| 3.8 | Generic construction of an oblivious fuzzy extractor (<i>RFE</i>) on top of an oblivious secure sketch scheme (<i>RSS</i>). | 98 |
| 3.9 | EER interval from the uncertainty of the FMR and FNMR. | 103 |
| 3.10 | Illustration of a Susceptible-Infectious (SI) model for the zombie attack. | 104 |

List of Tables

| | | |
|------|---|----|
| 1.1 | Comparison Between Personal Data and Biometric Data | 6 |
| 1.2 | Short comparison between Steganography and Cryptography | 13 |
| 2.1 | Comparison of cooling methods for our simulated annealing. | 28 |
| 2.2 | Summary of the experiments of the space partitioning algorithm. | 30 |
| 2.3 | Critical population for various threshold and template sizes for $\gamma = 0.01$ | 39 |
| 2.4 | Critical population comparison in function of FMR with $\lambda = 2$ | 41 |
| 2.5 | FMR maximal needed in function of the population size of the system for $\lambda = 2$ | 42 |
| 2.6 | Expected number of calls to <code>oracle</code> for the exhaustive search method Random Sampling with Replacement (RSR). | 46 |
| 2.7 | Summary of all leakage exploits and their complexities with α such that the occurrence of the rarest error is $n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$. The Distance-to-Threshold comparison determines if the leak occurs when $d(x, y) \leq \varepsilon$ (below) or when there is no distance requirement between x and y (both). For all the complexities, x and y are in \mathbb{Z}_q^n with $q \geq 2$ except for the minimal leakage where x and y are in \mathbb{Z}_2^n . The provided complexities represent worst-case scenarios, except for the accumulation attack where the result is the expectation. ¹ Note that the Big-O complexity of the optimal exhaustive search strategy, in the worst-case, is the same as the naive strategy as the minimum of $h(\cdot)$ is 0. | 47 |
| 2.8 | Summary of all leakage exploits and their complexities for the weaker attack model with α such that the occurrence of the rarest error is $n^{-\alpha}$ with $\alpha \in \mathbb{R}_{\geq 1}$. The Distance-to-Threshold comparison determines if the leak occurs when $d(x, y) \leq \varepsilon$ (below) or when there is no distance requirement between x and y (both). For all the complexities, x and y are in \mathbb{Z}_q^n with $q \geq 2$ except for the minimal leakage where x and y are in \mathbb{Z}_2^n . The provided complexities represent worst-case scenarios, except for the accumulation attack where the result is the expectation. ¹ Note that the Big-O complexity of the optimal exhaustive search strategy, in the worst-case, is the same as the naive strategy as the minimum of $h(\cdot)$ is 0. | 52 |
| 2.9 | Ratio between a κ -adaptive attacker and a 0-adaptive attacker in function of n , ε , N , a and κ for the lower bound. | 58 |
| 2.10 | Bounds for the number of operations for an outsider, in function of n , N and ε | 59 |
| 3.1 | Concrete security estimates for the cases $n = 128, 256$, $b = 2, 4, 8$, and $t = (b - 1)n/8$. In this table, C_{Guess} estimates the complexity of the guessing attack as $2^{\mu_{G^n}}$ while C_S estimates the complexity of the attack based on solving DLP, Knapsack, and SVP as $2^{0.2n}$ | 83 |
| 3.2 | Experimental results for Sketch and Rec implementations over the parameters $n = 640$, $b \in \{2, 4, 8\}$, and $t = (b - 1)n/8$. Timings have been averaged over 100 iterations. | 86 |

Résumé

Cette thèse vise à mettre en évidence les vulnérabilités des systèmes biométriques et à proposer des solutions pour renforcer la sécurité de ces données.

Un système biométrique permet d'authentifier ou d'identifier un individu en utilisant des caractéristiques physiques ou comportementales, telles que les empreintes digitales. Pour des raisons de sécurité, ces données ne sont pas utilisées en clair mais sont transformées en gabarits, rendant difficile la reconstitution des données originales. Cette transformation assure le respect de la vie privée des individus tout en permettant une authentification et une identification précises. En raison de leur utilisation analogue, nous avons étudié les données biométriques de manière similaire aux mots de passe en cryptographie. Plus précisément, nous avons d'abord étudié la probabilité d'occurrence d'une quasi-collision, c'est-à-dire la probabilité que deux gabarits de deux utilisateurs distincts soient proches. Les quasi-collisions posent problème car elles dégradent la capacité de reconnaissance du système et peuvent être exploitées par un attaquant cherchant à usurper l'identité de plusieurs utilisateurs. Pour éviter ces inconvénients, nous avons établi une borne sur la taille de la base de données pour prévenir les quasi-collisions et introduit un score pour aider à paramétrer les algorithmes de reconnaissance biométrique.

Ensuite, nous avons étudié les attaques par recherche exhaustive sur les données biométriques. Nous avons d'abord examiné les attaques ciblées, visant un utilisateur particulier, en étudiant la probabilité qu'un attaquant réussisse à usurper l'identité d'un utilisateur choisi dans différents scénarios. Cette étude nous a permis de définir des bornes de sécurité pour les bases de données de gabarits et de fournir des recommandations concernant les paramètres de sécurité pour les systèmes biométriques. Nous avons également investigué les attaques non-ciblées, où l'attaquant ne vise aucun utilisateur en particulier, pour évaluer la probabilité qu'un ou plusieurs attaquants réussissent à usurper l'identité de quelqu'un dans une base de données. Même si la probabilité d'usurper l'identité d'un individu spécifique est faible, il peut être facile d'usurper l'identité de quelqu'un lorsque la base de données est grande, de la même manière qu'il est probable que "0000" soit le mot de passe de quelqu'un dans une grande base de données. Cette analyse nous a permis de compléter notre investigation de la sécurité des données biométriques et de caractériser leurs limites.

Les attaques mentionnées ci-dessus s'appliquent principalement hors ligne. En ligne, ces attaques sont généralement détectées ou des contre-mesures sont mises en place pour ralentir les attaquants, comme l'augmentation du temps d'attente entre les tentatives. Pour limiter les problèmes liés aux attaques hors ligne, nous avons développé deux nouveaux protocoles d'authentification biométrique résistants aux attaques hors ligne. Le premier protocole utilise une preuve à divulgation nulle de connaissances pour garantir qu'un client malveillant, même avec des ressources de calcul illimitées, ne puisse obtenir aucune information utile du serveur pour effectuer une recherche exhaustive hors ligne. Le second protocole permet de corriger la donnée biométrique fournie par le client, conçu de telle sorte qu'un client malveillant avec une capacité de calcul polynomiale ne puisse obtenir aucune information utile.

Mots clés : sécurité biométrique ; transformations biométriques ; authentification biométrique ; identification biométrique ; extracteur flou réutilisable ; esquisse sécurisée réutilisable ; extracteur flou opaque ; esquisse sécurisée opaque ; calcul sécurisé ; preuve par jeu ; divulgation nulle de connaissances ; distance obfusquée ; correspondance floue ; distance de hamming ; fuite d'information ; problème du collecteur de coupons ; problème de la chaîne la plus proche ; quasi-collisions
