

# Trust the System

## Auditing Privacy-preserving Medical Data Analysis in a Distributed Manner

Jorrit van Assen



# Trust the System

## Auditing Privacy-preserving Medical Data Analysis in a Distributed Manner

by

Jorrit van Assen

to obtain the degree of Master of Science  
at the Delft University of Technology,  
to be defended publicly on Thursday September 28th, 2023 at 9:30 AM.

Student number: 5652693  
Project duration: January 16th, 2023 – September 28th, 2023  
Thesis committee: Dr. Z. Erkin, TU Delft, supervisor  
Dr. M. Khosla, TU Delft

An electronic version of this thesis is available at <http://repository.tudelft.nl/>.

# Abstract

Recent developments in the capability and availability of small internet of things devices has meant that networked medical devices, like networked implants and wearable monitors, have become more widespread. This data is invaluable for solving pressing global healthcare concerns, like effectively monitoring and treating heart patients. The European Union has announced plans to create an international collaborative network for sharing medical data. However, such a system will have to overcome some major unsolved issues regarding security and privacy. Citizens surveys have stressed the importance of privacy protection and transparency in recipients. Governments have appointed administrative bodies tasked with supervising the processing of personal data, or assuring healthcare quality. However, medical healthcare providers have signalled concern with unrestricted governmental access to patient data. In this thesis, we propose a system for auditable medical data sharing compatible with privacy-preserving technologies. We demonstrate a method to securely generate encryption keys which are recoverable using an audit key. We combine this with distributed key generation to create a board of trusted members, with each a share of the audit key. Board members can work together to collaboratively audit communication between healthcare providers and medical researchers. We demonstrate that the key generation is secure and efficient. We show that auditability is guaranteed under the assumptions that at least one of the communicating parties is honest. Our system bridges the gap between privacy-preserving medical data analysis and governing capabilities by assuring auditability without handing this power over to a single party. In real world scenarios, this system can be used to create international level of data sharing, as is explored for the European Health Data Space. The data inspection can be combined with already existing legislative power to detect fraudulent behavior and perform physical audits when required. The system can be extended to facilitate reproducible medical research.

# Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>                        | <b>1</b>  |
| 1.1      | Research Gap . . . . .                     | 4         |
| 1.2      | Contributions . . . . .                    | 4         |
| 1.3      | Outline . . . . .                          | 5         |
| <b>2</b> | <b>Preliminaries</b>                       | <b>6</b>  |
| 2.1      | Communication Channels . . . . .           | 6         |
| 2.2      | Hybrid encryption scheme . . . . .         | 6         |
| 2.2.1    | Fujisaki-Okamoto transformation . . . . .  | 7         |
| 2.2.2    | PSEC . . . . .                             | 7         |
| 2.2.3    | EC-IES . . . . .                           | 7         |
| 2.3      | Distributed Key Generation . . . . .       | 7         |
| 2.4      | Distributed Ledgers . . . . .              | 8         |
| <b>3</b> | <b>Related Work</b>                        | <b>9</b>  |
| 3.1      | Privacy preserving analysis . . . . .      | 9         |
| 3.2      | Medical Data Sharing . . . . .             | 10        |
| 3.3      | Auditing in Medical Data Sharing . . . . . | 10        |
| 3.3.1    | Auditing Through Trusted Parties . . . . . | 11        |
| 3.3.2    | Auditing in Zero Knowledge . . . . .       | 11        |
| <b>4</b> | <b>Auditable and Secure Data Sharing</b>   | <b>13</b> |
| 4.1      | Stakeholders . . . . .                     | 13        |
| 4.2      | Threat Model . . . . .                     | 14        |
| 4.3      | System Requirements . . . . .              | 15        |
| 4.4      | Fraud protection . . . . .                 | 16        |
| 4.5      | Notation . . . . .                         | 17        |
| 4.6      | Assumptions . . . . .                      | 17        |
| 4.7      | Our Proposal . . . . .                     | 17        |
| 4.7.1    | Public Ledger . . . . .                    | 17        |
| 4.7.2    | Interactions . . . . .                     | 18        |

|          |  |           |
|----------|--|-----------|
| <b>5</b> | <b>Evaluation</b>  | <b>25</b> |
| 5.1      | Complexity Analysis . . . . .                              | 25        |
| 5.1.1    | Computational Complexity . . . . .                         | 25        |
| 5.1.2    | Message Complexity . . . . .                               | 26        |
| 5.2      | Security Analysis . . . . .                                | 27        |
| 5.2.1    | Correctness of the Communication Keys Generation .         | 27        |
| 5.2.2    | Proof of Auditability . . . . .                            | 27        |
| 5.2.3    | Proof of Key Authenticity . . . . .                        | 28        |
| 5.2.4    | Security Sketch of Communication Keys . . . . .            | 29        |
| <b>6</b> | <b>Discussion</b>  | <b>30</b> |
| <b>7</b> | <b>Future Work</b>   | <b>33</b> |
| 7.1      | Privacy-preserving Outsourcing of Data Analytics . . . . . | 33        |
| 7.2      | Reproducible Medical Research . . . . .                    | 34        |
| 7.3      | Further Auditable Data Sharing . . . . .                   | 34        |
| <b>8</b> | <b>Conclusion</b>  | <b>35</b> |

# Preface

This thesis marks the conclusion of my master's in Computer Science at the TU Delft. I will be forever grateful for the opportunities that my study has given me. I would like to express my sincere gratitude to the people who have helped me with this work over the last 9 months. Foremost, this thesis would not have been possible without the help and continuous support of my supervisor Dr. Zekeriya Erkin. Many thanks to Florine Dekkers, Tianyu Li, Jelle Vos and Roland Kromes for providing me with feedback and new insights. I am also grateful for the tips and mental support of my fellow Cyber Security students. Special thanks go to Ivo and Andrei, with whom I have shared the majority of my thesis experience. Lastly, I would like to thank my family and friends for their dedicated and unwavering support.

*Jorrit van Assen  
Delft, September 2023*



*ich möchte gern zwei kleine Hunde sein  
und miteinander spielen.*

—Friedrich Torberg, *Ballade der großen Müdigkeit*

# Chapter 1

## Introduction

The rapid developments in digital technologies are double-edged sword for hospitals. They provide great opportunity with smart medical devices and automated quality control tools on one hand, but also with greater attack surface and digital complexities on the other. In [1], Scheibner et al. noted a “biomedical research shifted from paradigm has been characterized by a shift from intrainstitutional research toward multiple collaborating institutions operating at an interinstitutional, national or international level for multisite research projects”. How to facilitate this shift in a secure and privacy-preserving manner, we adres in this thesis.

**Digital struggles.** Public healthcare has been unable to keep up with the rapid technological developments in the digital space. An important aspect of healthcare is certification, and the lengthy bureaucratic procedures associated with certifying medical equipment means that most electronic devices are already years old before they become available to hospitals. For digital devices, the lifespan of hard- and software is often a couple of years. Although medical devices are purposely designed for medical applications, and in general designed with this delay in mind, often times they are based on software and hardware modules which are not. This means that medical apparatus is often times harder to secure against cyber threats, especially when these devices are connected to the internet. An example of this was during the WannaCry ransomware attacks in April and May 2017. Up to 70000 devices of the National Hospital Service of the United Kingdom were compromised. The malware made use of a security vulnerability that had been patched the month before in supported versions of the Windows operating system. However, devices running Windows XP or which were not updated were still vulnerable. The attack had a significant effect on the continuity of healthcare, and some hospitals had to cancel non-critical operations[2]. Better awareness of cybersecurity in medical facilities is required, especially when more patient information is collected and shared. However,

more importantly, these events signify the need for a different approach; A system designed to failing well.

**Making healthcare future-proof** The graying of European population and the increasing public expenditure on healthcare calls for a better approach. In the joint priorities of the EU institutions for 2021 until 2024, healthcare is an import part of “Promoting our European way of life”. It calls for better responses to international health threats, more affordable and accessible healthcare, beating cancer and a so-called European Data Space. In a recent impact assessment report on the European Data Space gave the following definition for the European Data Space.

The general objective of the intervention is to establish a genuine single market for digital health and to ensure that individuals have access to and control over their own health data, can benefit from a wealth of innovative health products and services based on health data use and reuse, and that researchers, innovators, policy-makers and regulators can make the most of the available health data for their work, while preserving trust and security.

The usefulness of European Health Data Space (EHDS) is twofold. The primary use of the European Data Space is to improve the current healthcare system. The patient stays in control of their electronic health records, and different healthcare providers can make use of the same patient dossier. Secondly, this medical data could also be used for research and policymaking purposes. This secondary purpose would help researchers, innovators, and regulators to base their work on actual data and improve healthcare practices for the general population. Both of these use cases requires strongly security guarantees and trust from the patients.

**Potential fall blocks.** An article in the Dutch newspaper NRC clearly demonstrated how conflicts in privacy guarantees for medical data can create obstacles for important research. To optimize funding for mental healthcare, the NZa (Dutch Care Authority) wants to predict what types of healthcare are needed in the future. To this end, it requires practitioners in mental and forensic healthcare to fill in a form for each patient treated. This form includes questions like: “Does this patient show destructive behavior?”, “Is this patient aggressive?”, “Does the patient abuse drugs?”, “Repeated self-harm?” and “Have they attempted suicide?” [3]. Although personal identifiable information, like date of birth or name, will not be included in the reports, it is possible to re-identify patients by combining the form with other data sources. Critics of this approach have united under the name “vertrouwenindeggz” with the goal to sue the NZa. Their claim is that forcing healthcare providers to fill in this form is an “unacceptable violation of the medical confidentiality and client privacy” [4].



However, it is likely that such detailed records of individual patients are not required for the NZa to make its estimations. It could instead choose to query information in a more limited scope. For example, a policymaker that has a legitimate interest in knowing the amount of terminally ill patients in South-Holland does not need to know the names of these patients and their respective medical history. Preferably, it should not even know the exact number of terminally ill patients per hospital. Just one final number that gives a direct response to their initial query: The number of terminally ill patients in South-Holland. Depending on the use case, a small error on this number might not be significant. A well-defined query makes sure that its use is limited, the ground testable and the potential fallout, in case of data leakage, manageable. Of course, hidden in this query are many complexities; A usable definition of a terminally ill patient needs to be defined, the right data providers identified, etc. Additionally, the policymaker might want to include more conditional statements or query multiple times. Still, queries on medical health data should provide the exact amount of data requested and not leak any additional information. In this way, it is possible to perform useful analysis, without violating the privacy of patients.

These observations are also reflected in a report by TEHDAS. TEHDAS is an EU-sponsored joint action that explores the steps towards EHDS. Its main focus is to develop European principles for the secondary use of health data. Their 2022 survey incorporated the opinion of 5932 citizens, primarily from France, Belgium, and the United Kingdom[5]. A majority of the respondents deemed secondary use of data for the common good as a valid purpose. A recurring theme throughout all responses was that citizens perceive health data as a piece of them, should benefit them, and therefore they should retain some forms of control over them. This holds true even when anonymized, although the researchers found that the respondents were more open to a broader use of their data. The citizens expressed fear of being re-identified and potential harmful consequences, like discrimination. They argued in favor for safe IT solutions, diversity of views and backgrounds of stakeholders, a framework for accountability and transparent intentions of data users.

Interestingly enough, stakeholders of EHDS have argued in favor of more responsibility for citizens themselves, with citizens responsible for managing their data and being aware when to opt out of research. This is also reflected in the numerous previous publications on EH systems that employ different techniques to give citizens digital control over their patient data [6][7]. However, relying just on informed consent is not an adequate solution. Data sharing and analytics is a very abstract and technical topic. Asking people to understand the process is a heavy burden, disproportional impact of their time. A similar legal document, the terms of service are often ignored because of their length and difficulty [8] [9]. Additionally, creating technical solutions for citizens to control their data usage risks the exclusion of less

technical versed individuals. Therefore, consent forms should be clear and concise, allowing individuals to express their values and should not rely on complicated technical solutions. General ethical standards could be ensured by a governing body that uses the transparency of the data analysis process to this end.

## 1.1 Research Gap

In recent years, multiple directions for improving digital healthcare have been explored. Privacy-preserving data analysis techniques have been proposed which are deployable in real-time. Additionally, extensive research has been done into utilizing distributed ledgers for tracking medical data sharing. Ensuring auditing in these systems, however, is still not completely solved. Existing solutions put limitations on data types, have significant computational overhead, or utilize trusted third parties.

Combining this, a solution for a medical data sharing platform for the purpose of research and analysis should be general purpose, secure, transparent and safeguard the patient privacy, without trusted third parties. To the best of our knowledge, there are no current solutions which manage to combine these three aspects. A suitable protocol must allow for large scale medical research to be performed, while the privacy loss of patients should be specifiable. The method should be compatible with data gathered from Networked Medical Devices and should provide minimal computational cost for the healthcare providers. This leads us to the following research question:

“How can analytics on medical data be performed in a secure and privacy-preserving manner, without a trusted third party and while keeping the ability for an auditor to test the soundness of the procedure?”

## 1.2 Contributions

In this thesis, we provide the following contributions.

1. We define the requirements and threat model for privacy-preserving medical data analysis that incorporated auditability for a consortium of auditors.
2. We propose a generic key generation protocol, which enables third parties to recover keys under the condition that at least one of the two generators is non-malicious.
3. We design a system based on our key generation protocol and distributed ledgers to achieve auditable and privacy-preserving data sharing. We perform a complexity analysis and security analysis of the protocol.

### 1.3 Outline

The rest of the thesis is organized as follows. Chapter 2 covers the preliminaries required for the protocol. Chapter 3 gives an overview of work related to privacy-preserving medical data analysis and auditable systems. Chapter 4 gives the threat model, requirements of the system, and formal description of the protocol. Chapter 5 provides the complexity, and security analysis of the protocol. Chapter 6 discusses the results, and reflects on the feasibility of such systems in the current healthcare landscape. And finally, Chapter 8 concludes this thesis.

## Chapter 2

# Preliminaries

In this chapter, we will cryptographic concepts used in our thesis. We assume the reader familiar with basic concepts of cryptography and information theory. More specifically, public key encryption, symmetric key encryption, Readers unfamiliar with these concepts, we refer to “Introduction to Modern Cryptography“ by Jonathan Katz and Yehuda Lindell<sup>1</sup> and “A Graduate Course in Applied Cryptography” by Dan Boneh and Victor Shoup<sup>2</sup>.

### 2.1 Communication Channels

To evaluate the security of protocols involving multiple parties, the security of their communications must be included. We will differentiate two type of communication channels: Authenticated channels, and secure channels.

**Authenticated channel.** An authenticated channel is a communication channel between two parties where messages are authenticated. This means that the identity of the sender is guaranteed and that the integrity of the messages are assured. However, the content of the messages send over the communication channel are visible to observers.

**Secure channel.** Secure channels have the same properties as authenticated channels. However, they contain as additional guarantee that messages are undecipherable to observers. This can be achieved using authenticated encryption.

### 2.2 Hybrid encryption scheme

Hybrid encryption scheme combine symmetric and asymmetric cryptography into a single algorithm. Asymmetric cryptography has the advantage of allowing parties to encrypt messages using public keys. However, they are

---

<sup>1</sup><https://www.cs.umd.edu/~jkatz/imc.html>

<sup>2</sup><https://toc.cryptobook.us/>

more expensive in terms of computational and space complexity compared to symmetric cryptographic methods. By combining the two methods, the strength of both systems can be exploited.

In hybrid encryption schemes, the sender generates a new key used for encrypting the message. The generated key is encapsulated using public-key cryptography. The encapsulated key and encrypted message are sent to the receiver. The receiver can recover the symmetric-key from the encapsulation using their private key. The symmetric-key is then used for decrypting the ciphertext to recover the original message.

### 2.2.1 Fujisaki-Okamoto transformation

The Fujisaki-Okamoto transformation is a general way to convert any weak symmetric and weak asymmetric encryption schemes into a IND-CCA secure encryption schemes. The transformation was first published in 1999 [10] and later revised in 2013 [11]. Based on the transformations, the authors proposed PSEC. It is possible to use Fujisaki-Okamoto transformation without using a symmetric encryption scheme.

### 2.2.2 PSEC

PSEC stands for Provably Secure Elliptic Curve Encryption Scheme. PSEC-KEM is a public encryption scheme that uses the elliptic curve ElGamal trapdoor function, two hash function and a symmetric key encryption scheme. [12] The scheme is IND-CCA secure in the random oracle model under the EC-DH assumption. Two other versions, PSEC-1 and PSEC-3, have existed which have been withdrawn from the CRYPTREC 2000 project. [13] PSEC-KEM can be used as an authenticated encryption protocol.

### 2.2.3 EC-IES

The Elliptic Curve integrated Encryption Scheme (EC-IES) is a public encryption scheme that uses Diffie-Hellman key exchange to directly compute a key used for symmetric encryption. [14] The scheme is IND-CCA secure in the random oracle model under the EC-CDH assumption. There exists multiple standards of EC-IES. [15]

## 2.3 Distributed Key Generation

The concept of Distributed Key Generation (DKG) was introduced by Pedersen in “A Threshold Cryptosystem without a Trusted Party” [16]. It allows the oblivious generation of a shared secret by members of a group, such that a group of  $k$  members are required to reconstruct the secret. Whereas Shamir Secret Sharing requires a trusted third party as an oracle for generating the

original secret, DKG can be performed without the need for a trusted third party. Instead, a concept called verifiable secret sharing (VSS) is employed.

VSS allows a secret to be divided into multiple pieces, such that combining the shares yields the original secret. Every receiver of the secret can verify that the secret is a valid piece of the original secret. VSS allows the distributed generation of a key, even with the presence of malicious actors. However, to ensure the security properties of DKG, the number of malicious actors must be bounded by  $l$  where  $l \geq 2k - 1$ .

## 2.4 Distributed Ledgers

A distributed ledger is a history of data transactions shared by multiple parties. The order of the history is agreed upon, and each party has a copy of the transactions. New transactions are added using a consensus protocol.

Blockchain is a form of a distributed ledger, where the transactions, called blocks, receive a timestamp and a reference to the previous block. Blockchain was originally proposed as a distributed ledger to track the cryptocurrency Bitcoin. However, the use of blockchain outside of cryptocurrencies has been extensively researched.

## Chapter 3

# Related Work

In this chapter, we will discuss works that are related to privacy-preserving medical data analytics and auditable data sharing.

### 3.1 Privacy preserving analysis

Privacy-preserving analysis allows data processing in encrypted form. This has the benefit that auditing the processors will not leak information about the underlying sensitive data. Often these works employ either multi-party computation or homomorphic calculations. These works are interesting for this thesis since they often involve a similar scenario to ours.

In [17], Blanton et al. present a general and scalable strategy to enable privacy-preserving learning of generative and discriminative machine learning. Participants are divided into data owners, computational service providers and output recipients. One participant may be part of multiple groups without damaging the integrity of the protocol. By using secure multi-party computation, can protect the security of the data, however this introduces overhead in the computations. This strategy is demonstrated by performing the variable assignment problem, without direct access to the medical records. Speed-ups over a naive implementation are made by eliminating divisions and precomputing logarithms.

In [18], the authors propose a system framework that allows to delegate the inference process on medical data to a computation service provider. The medical diagnostics is completely done in the ciphertext domain, and thus the service provider will have no information on the data processed or the result of the inference. The system framework solely relies on secret sharing instead of homomorphic encryption and garbled circuits in order to reduce time and message complexity. The protocol is split into a preprocessing and an online phase. A large part of the calculations are done during the preprocessing phase to limit the number of interactions.

## 3.2 Medical Data Sharing

MeDShare is a system proposed by Xia et al. [19] that enables the processing and sharing of medical data in a trustless environment. The system aims to solve two issues. Entities performing actions on data without knowledge of the data owner, compromising privacy and value of data. Entities accessing processing results with different intended recipients. The proposed systems provide identity management, decentralized and centralized access, data access revocation and tamperproof data auditing. This is achieved by separating 4 layers: The User layer, Data Query layer, Data Structuring and Provenance layer and the Database layer. Entities interface with the User layer using predefined query and response structures. The Data Query layer acts as a transformer between this structure and environments outside the system. The Data Structuring and Provenance layer is responsible for processing requests for data from the Database layer. Lastly, the Database layer consists of data providers with established database infrastructure which are connected to the system.

## 3.3 Auditing in Medical Data Sharing

Bonyuet presented in [20] the impact of blockchain in the audit profession. It discusses both ways in which blockchain can aid in making financial systems auditable, and risks associated with blockchain. Prior research by Fanning & Centers [21] had established that blockchain might have impact on auditing. Because of the design of blockchain, it can represent a triple entry accounting system with immutable transactions, time stamped and encrypted [22]. However, blockchain will not guarantee all transactions to be sound. The transactions could still be unauthorized, executed between related parties, influenced by an “off-chain” agreement or incorrectly classified. Therefore, auditors will still be required to perform audits on the soundness of the transactions. However, blockchain can aid in speeding up the auditing process. “The focus of an audit will shift from record tracing and verification to complex analysis, such as systemic evaluation, risk assessment, predictive audits, and fraud detection.” The adoption of blockchain can speed up the auditing process and therefore increase the level of auditability in general. Additionally, smart contract might be useable for automating manual auditing tasks. The authors also point out the risks of using blockchain for auditing. Traditional blockchains, like bitcoin, use proof of work as a consensus technique, which is susceptible to a 51% attack. Additionally, the authors discuss tradeoffs between public and private blockchains. Although public blockchains avoid centralization of control, it comes with the drawback of enabling unauthorized parties from reading or even writing to this Blockchain. For private blockchains, or permissioned blockchains, the au-



thors claim that since a single party has control over the validation, they would then be able to modify the transaction history as needed.

However, Yannis Bakos et al. argue in [23] that permissioned blockchain can actually increase decentralization when the blockchain is designed carefully. They observe that in order to ensure quality control and coordination of system development and evolution of open access and fully distributed environments, *de facto* centralization arises naturally. The authors name expertise, reputation, time, or money as potential obstacles that can limit the participants to a selected few. The higher these costs are, the fewer the people that want to participate, which contributes to this centralization in practice. This also translates to permissionless blockchains. Permissionless blockchains like Bitcoins often lack preventive measures for countering centralization. As a result, centralization emerges in practice, for example through large mining pools who provide most of the work required for consensus.

### 3.3.1 Auditing Through Trusted Parties

Chen et al. describe a secure storage and sharing system for medical records using blockchain [24]. The author aims to streamline medical record sharing and solve drawbacks of previous work like reliance on trusted third parties and lack of patient control. The authors chose to avoid high storage cost by only tracking metadata on the blockchain, the medical data itself is stored on a cloud instance. “The blockchain is viewed as a storage supply chain in which every operation may be verified, accountable and immutable. Such inherent characteristics make it a potential solution for healthcare data systems that concerns both sharing and patient privacy.” By encrypting all data using a public key of the patient, the patient has full control over who can read the data. Privacy is preserved through techniques like hiding the identity of the patient by using different public keys for different hospitals.

### 3.3.2 Auditing in Zero Knowledge

Narula et al. propose zkLedger [25] in 2018, which enables auditing of assets in a ledger without revealing private information. The authors identify that distributed ledger either expose sensitive information publicly, or permits third-party auditors to decrypt sensitive content. zkLedger enables parties to create digital asset transactions which are visible only to those who are involved, but verifiable by everyone who is not. This is achieved using Schnorr non-interactive zero-knowledge proofs. zkLedger guarantees fast and provably correct auditing. Additionally, zkLedger guarantees completeness, preventing parties from hiding transactions from an auditor. The system is limited to numerical systems, and is specifically designed for tracking financial systems. Auditors are able to determine sums using only one round

of communication. More complex queries, however, require more involved communications with the audited party. The authors demonstrate how more complex queries like determining moving averages, variance, standard deviation, and ratios can be determined. Caching is employed to significantly reduce the cost of auditing, however, even without caching auditing stays viable. 100K transactions can be audited in 3.5 seconds, without cached commitments. time complexity of auditing cached transaction commitments is constant, whereas time complexity of un-cached transaction commitments grows linearly with the number of transactions in the ledger. Additionally, the time complexity of verifying transactions increases quadratically with the number of banks.

Zhao et al. introduce in [26] a new controllably linkable group signature scheme CL-GS and build on this a blockchain-based auditable privacy-preserving data classification (PPDC) scheme for IoT. PPDC allows an auditor to check the behavior of a semihonest data process and data center. IoT devices collect data and send this signed to the data center. Data processors interact with the data center to execute privacy-preserving data classification. The auditor can verify the correctness of the data classification results and records the auditing results on the blockchain. The basis of the system are controllably linkable group signatures (CL-GS), which requires a link key to check if two group signatures come from the same signer.

Zheng et al. propose in [27] a simplification of the PPDC scheme that reduces the trust requirements on the auditor. Since the auditing phase of the original PPDC scheme, all participants are perform multiple rounds of interactions the user cost increased and practicability is lowered. In the proposed scheme called VeriDC, verifiable proofs are immediately generated, and the data center can verify the classification results without an auditor.

zkrcChain by Xu et al. [28] focuses on auditing data as opposed to digital assets. Identities are kept public, but the data auditing can be performed by multiple provers, who collaboratively generate a joint range proof for one verifier. The system is built on the hyper-ledger fabric consortium blockchain and multi-party proofs based on Bulletproofs, zero-knowledge range proofs. It allows for arbitrary-range proof generation and verification and batch verification. The verifying party defines the upper and lower limit of the range and adds this information to the smart contract. The provers read the private data from the off-chain database and generate the proof. A dealer peer aggregates the proofs and which uploads the single proof to the ledger. Afterward, the verifying party can execute the smart contract to verify the proof.

## Chapter 4

# Auditable and Secure Data Sharing

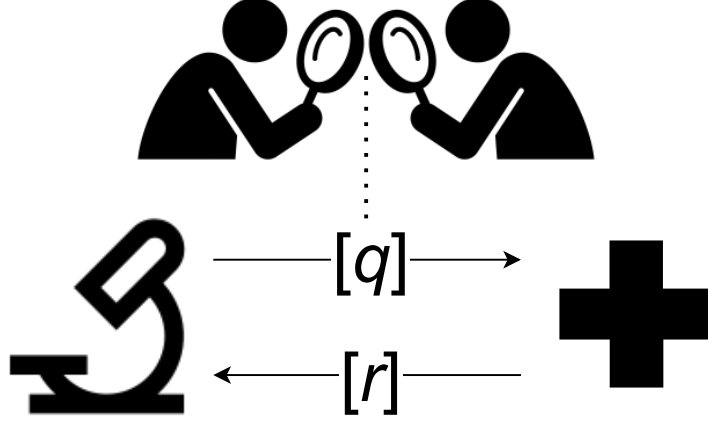
In this chapter, we present the design of our protocol for Auditable and Privacy-Preserving Medical Data Sharing. The goal of the protocol is to facilitate an ecosystem where medical data can be shared in a privacy-preserving way while retaining reproducibility and accountability. Research groups and hospitals join a network supervised by a board of trusted entities. Research groups send queries to Hospitals with the aim to create useful insights based on sensitive medical data. Hospitals execute the queries and send back the results to the querent. The protocol defines a systematic way for sending and securing these query-response communications. The messages are encrypted using special keys which are generated by the involved research group and hospital. The generation process allows the board to recover these keys, however this requires a threshold of board members involved to be met. All communications are stored in a distributed ledger, allowing the messages to be audited by the board.

In section 4.1, we give definitions for the stakeholders. In section 4.2 we establish our threat model. We follow this up in section 4.3 with our definition of reproducibility and accountability and the system requirements of the protocol. Section 4.5 gives an overview of the notation used in this thesis. Lastly, section 4.7 gives a detailed description of the proposed protocol.

### 4.1 Stakeholders

We identify three types of parties for this protocol: Hospitals ( $\mathcal{H}$ ), Research Groups ( $\mathcal{R}$ ) and Board Members ( $\mathcal{M}$ ). An organization may be multiple party types at once. Figure 4.1 provides an overview of the setting.

**Hospitals.** We use hospitals as an umbrella term for any organization that has access to patient data and executes queries. We denote an instance of a hospital with the symbol  $\mathcal{H}$ . Hospitals have a collection of patient data

Figure 4.1: Overview of the  $\mathcal{H}$ ,  $\mathcal{R}$  and  $\mathcal{M}$  party types.

with permission from patients to purpose their medical data for secondary usage. We assume that the hospital retains full control over this data and shares the data solely through the described protocol.

**Research group.** A research group denotes any type of instance, public or private, which requires medical data for research or decision-making purposes. We denote an instance of a research group with  $\mathcal{R}$ . Research groups create queries which hospitals evaluate. To ensure that research groups do not lose their competitive edge, queries and query-results remain confidential between the research group and hospital until audited by a delegation of the board.

**Board member.** A board member is a designated instance that is entrusted with a key share required to audit the system. They may be bodies appointed by governments, hospitals, or journal boards. Together with other board members, the board forms a consortium that monitors queries and their results. Given enough members interested in verifying specific data exchanges, the contents of the queries and results can be decrypted and audited.

## 4.2 Threat Model

We assume the hospital to be a semi-honest adversary. In related work, systems have been proposed where the patient remained in control of their data. However, this comes with a large increase of computational complex-

ity, storage increase and communication overhead. Additionally, hospital devices that are used in testing the patients will be operated and controlled by hospitals. This means that the patient has to trust the hospital to not have copied the information while it was created. In general, hospitals are assumed to be non-malicious. There are security measures in place that aid with ensuring the integrity of operations and preventing, for example, individuals to gain access to patient information directly. Although it would allow for stronger protocols, including Hospitals in the threat model would create a more significantly more complex system, which could increase the cost of deployment in the real world and would prevent hospitals from gradually increasing their security systems.

On the other hand, we assume some or all research groups to be a malicious adversary. The protocol should retain trustworthiness even when research groups actively try to increase the system. In this way, valid transcripts of data transfers can prove the legitimacy of research performed and aid in the repeatability and transparency of research. Research groups may actively try to gain information shared in data transfers where they are not the intended recipient.

We assume that less than  $t$  board members can be corrupted during any run of the distributed key protocol. In this way, the confidentiality of the queries and the privacy of the patients are ensured.

### 4.3 System Requirements

We need to define reproducibility and accountability first. We explain reproducibility in terms of the need for repeatability in the system. If any party has access to the full communication transcript, it would be able to repeat the exact same steps as denoted and achieve the exact same results. the protocol's For accountability, we need to ensure that parties take responsibility for their communications. To achieve the reproducibility and accountability, we identify the following system requirements.

**Confidentiality.** All key commitments, queries, and responses sent over the network are encrypted in a way such that only participants with access to the private key are able to recover the underlying message.

**Auditing.** To be able to audit encrypted queries and responses sent over the network by  $\mathcal{H}$  and  $\mathcal{R}$ , a consortium of board members should be able to reconstruct the communication keys used for encryption.

**Definition 1** *For all encrypted query-response pairs  $\langle c_q, c_r \rangle = \langle \text{enc}_{k_q}(q), \text{enc}_{k_r}(r) \rangle$ , with respective key commitment  $\langle c_{k_q}, c_{k_r} \rangle$ , the system is auditable if and only*

if, at least  $t$  board members can recover  $k'_q, k'_r$  such that  $\text{dec}_{k'_q}(c_q) = q$  and  $\text{dec}_{k'_r}(c_r)$ .

In order to audit the transactions, a consortium of  $t$  Board Members can collaboratively decrypt the contents of a ledger-entries. The protocol defines a systematic way for securing these query-response communications. The security is done in a way such that for a specific chain of Hospital and research group communication. Board Members private key shares can be revoked, and new members can be added to the board. In the case that new members are added to the board, they will only be able to audit transactions that were performed after admission.

**Integrity.** The ledger keeps track of all queries and results. When two parties agree on a transaction, it must remain integrous. A malicious research group must not be able to change the query transaction once it has been committed to the ledger. Additionally, a malicious research group must also not be able to fake data results transactions.

**Non-repudiation.** If a hospital has evaluated a query and transferred the result to the corresponding research group, it is unable to remove the transaction from the ledger or change the contents. It cannot repudiate the validity of the transaction, and given the same query and the same data entries it must be able to give the same output.

## 4.4 Fraud protection

When the system adheres to the requirements identified in Section 4.3 it protects against the following forms of fraud.

**Protection against cherry-picking of data.** The protocol protects against the cherry-picking of research data. Given a malicious research group that tries to force the results of their research, they could try to change their queries and targeted data source until a desired result is achieved. Their queries and data could then be linked to the research result, and seemingly valid research has been performed. However, given our assumptions, the board would see these disconnected request-response pairs on the ledger. These patterns are suspicious and can easily be recognized by the board and decrypted for inspection. The query information stored in the ledger is enough to decide whether cherry-picking has occurred, and the board does not need to request sensitive information from Hospitals.

**Protection against falsified results.** The protocol protects against research groups fabricating results. When a RG tracks its queries and publi-

cations using the protocol, it commits to the validity of the data. All results that are publicized are contained in the ledger. By directly linking a publication with the data sources used, it increases the validity of the research. When results in publications are doubted, it is possible to inspect queries made by the research party directly. This gives complete transparency in the basis of the results and the methodology used.

**Protection against fabrication of patient information.** The protocol does not directly guarantee the authenticity of the original data of the hospital. If the hospital data is corrupted and used for research, there is no way to detect this. In the case that the hospital is subject to inspection and problems are found within their databases, queries can be replayed to verify whether they were based on incorrect data. Publications based on results of queries on incorrect data can be investigated and possibly flagged.

## 4.5 Notation

Let  $\mathbb{G} := \{g^a : a = 0, \dots, q - 1\}$  denote a group, such that  $\mathbb{G}$  is a subset of  $\mathbb{Z}_p^*$  of cardinality  $q$ , where  $p$  and  $q$  are prime.

Table 4.1 provides an overview of all mathematical symbols used in this thesis.

## 4.6 Assumptions

We assume that a PKI system has been established, allowing all participants to have access to a secure channel for each other participant. Each participant can also verify signatures of every other participant.

## 4.7 Our Proposal

### 4.7.1 Public Ledger

To ensure the integrity and non-repudiation requirements of queries and responses, a permissioned ledger is utilized. The board maintains the public ledger and authorizes  $\mathcal{H}$  and  $\mathcal{R}$  willing to join the ledger. Parties are assigned a unique identifier,  $id_p$ , that is used in all further communications. The ledger serves as a method of ensuring consensus on a single shared history of communications. All authorized participants of the ledger can broadcast their communications to all  $\mathcal{M}$ , who perform consensus on the resulting history. Afterward, the accepted blocks are shared with all members of the ledger.

Table 4.1: Overview of symbols used.

| Symbol               | Definition                    |
|----------------------|-------------------------------|
| $q$                  | query                         |
| $r$                  | result                        |
| $pk_s$               | list of primary keys          |
| $g$                  | generator                     |
| $pk_p$               | public key of party $p$       |
| $sk_p$               | private key of party $p$      |
| $k$                  | symmetric key                 |
| $s_p$                | secret share of party $p$     |
| $c$                  | ciphertext                    |
| $\sigma$             | signature                     |
| $id_p$               | identifier of party $p$       |
| $\xleftarrow{R}$     | uniformly sample              |
| $\text{enc}_{pk}()$  | encrypt using $pk$            |
| $\text{dec}_{sk}()$  | decrypt using $sk$            |
| $\text{sign}_{sk}()$ | sign using $sk$               |
| $\text{hash}()$      | perform secure hash function  |
| $\text{hkdf}()$      | key derivation function       |
| $\text{eval}()$      | evaluate a query              |
| $\text{dp}()$        | perform differential privacy  |
| $\text{pub}()$       | publicize block on the ledger |

The general structure of blocks is shown in Table 4.2. A 256-bit integer is used for uniquely identifying blocks. Additionally, a single byte is used for denoting the type of block. The type is either `QUERY_KEY`, `RESPONSE_KEY`, `QUERY` or `RESPONSE`. These blocks extend upon the general structure with type specific values. Related communications are linked on the ledger; responses point to queries, and `QUERY_KEY` and `RESPONSE_KEY` blocks are always referenced by `QUERY` and `RESPONSE` respectively.

Table 4.2: The general structure of blocks on the ledger.

| symbol    | description              | size (bits) |
|-----------|--------------------------|-------------|
| $p_{own}$ | block identifier         | 256         |
| $type$    | type of the block        | 8           |
| ...       | (type specific contents) | ...         |

#### 4.7.2 Interactions

**Setup of the board members.** To allow for threshold decryption, all  $\mathcal{M}$  cooperate in a distributed key generation algorithm. For our proposal, we



Table 4.3: The structure of the QUERY\_KEY and RESPONSE\_KEY blocks.

| symbol             | description                 | size (bits) |
|--------------------|-----------------------------|-------------|
| $id_{\mathcal{H}}$ | identifier of $\mathcal{H}$ | 256         |
| $id_{\mathcal{R}}$ | identifier of $\mathcal{R}$ | 256         |
| $c$                | encrypted communication key | 512         |

Table 4.4: The structure of the QUERY block.

| symbol             | description                              | size (bits)   |
|--------------------|--|---------------|
| $p_{key}$          | points to the associated QUERY_KEY block | 256           |
| $id_{\mathcal{H}}$ | identifier of $\mathcal{H}$              | 256           |
| $id_{\mathcal{R}}$ | identifier of $\mathcal{R}$              | 256           |
| $enc_{k_q}(q)$     | encrypted query                          | variable size |

employ the key generation protocol by Kate et al.[29] A high level overview is provided in Algorithm 1. The protocol guarantees agreement in an asynchronous network of  $n \geq 3t + 2f + 1$  nodes, with at most  $t$  Byzantine adversaries and  $f$  crashes in a round. Further details can be found in the original paper. After completion of the algorithm, every board member  $\mathcal{M}$  has the public key and a share  $s_{\mathcal{M}}$ . The public key broadcasted among all members of the ledger to be used in further communication. The subroutine of the algorithm is performed in intervals to generate new shares for the same  $pk_{\mathcal{B}}$ .

**Generation of communication keys.** To allow  $\mathcal{H}$  and  $\mathcal{R}$  to communicate queries and results in a secure and auditable way, communications keys should be constructed in a way such that the board can recover the keys reliably. This is achieved in the following way. First,  $\mathcal{H}$  and  $\mathcal{R}$  collaboratively generate two keys, one keys for queries and one for results. Secondly, the keys are signed, encrypted and published on the ledger. Lastly, both  $\mathcal{H}$  and  $\mathcal{R}$  verify the published ciphertexts and abort if they detect tampering. The algorithm is shown in Algorithm 2. We now describe the algorithm in more detail.

Table 4.5: The structure of the RESPONSE block.

| symbol             | description                                 | size (bits)   |
|--------------------|---|---------------|
| $p_{key}$          | points to the associated RESPONSE_KEY block | 256           |
| $id_{\mathcal{H}}$ | identifier of $\mathcal{H}$                 | 256           |
| $id_{\mathcal{R}}$ | identifier of $\mathcal{R}$                 | 256           |
| $enc_{k_r}(h)$     | encrypted hash of all primary keys          | 256           |
| $enc_{k_r}(r')$    | encrypted response                          | variable size |

Table 4.6: The structure of the ABORT block.

| symbol             | description   | size (bits) |
|--------------------|---|-------------|
| $p_{key}$          | points to the associated RESPONSE_KEY/QUERY_KEY block | 256         |
| $id_{\mathcal{H}}$ | identifier of $\mathcal{H}$                           | 256         |
| $id_{\mathcal{R}}$ | identifier of $\mathcal{R}$                           | 256         |

**Algorithm 1** Setup of the board keys.**Input**

- $n$  number of board members
- $t$  reconstruction threshold
- $f$  number of crash board members allowed
- $i$  index of the board member
- $\tau$  session

**Output**

- $pk_B$  public key of the board
- $s_i$  secret share of board member  $i$

**Step 1: Generation of the keys.** For the creation of keys, we use the elliptic curve variant of the Extended Triple Diffie-Hellman key exchanges [30]. Ephemeral keys  $X_q, X_r, X_s, Y_q, Y_r,$  and  $Y_s$  are generated using random points and are exchanged between  $\mathcal{H}$  and  $\mathcal{R}$ . The ephemeral keys and long term public keys are used as key material for the HMAC-based Extract-and-Expand Key Derivation Function. Two keys are generated,  $k_q$  for encrypting queries and  $k_r$  for encrypting results. The advantage of using separate keys is that the board can choose whether to audit both queries and results, or only a single communication type. The keys are used as long term keys between  $\mathcal{H}$  and  $\mathcal{R}$ . Additionally,  $s$  is constructed as a shared secret seed.

**Step 2: Publication of the keys.** The key generation ensures that  $\mathcal{R}$  and  $\mathcal{H}$  have securely gained access to two common communication keys. However, for auditing, the board should be able to recover the keys. To this end,  $k_q$  and  $k_r$  are shared on the ledger protected under the public key of the board. The responsibility of publishing keys is divided among the two participants;  $k_q$  and  $k_r$  are published by  $\mathcal{R}$  and  $\mathcal{H}$  respectively. When the publication is divided among the two parties, it is vital to assure the authenticity of the published keys. If an incorrect key is published and not corrected, all future communication becomes inauditable. Therefore, we introduce two methods to attest the authenticity of the published keys. The first method allows the board to verify that both parties agree on the published key; The second method allows the other participant in the key generation to assure the correct key is published.

Given that the key  $k'_r$  published by  $\mathcal{H}$ ,  $\mathcal{H}$  needs to include proof that  $k'_r$

is indeed the generated key  $k_r$ . Since both the board,  $\mathcal{H}$  and  $\mathcal{R}$  are allowed to know the key, a signature by  $\mathcal{R}$  suffices. After key generation,  $\mathcal{R}$  uses EC-DSA to sign  $k_r$  using its private key  $sk_{\mathcal{R}}$ . The signature,  $\sigma_r$ , is sent to  $\mathcal{H}$  over an authenticated channel and  $\mathcal{H}$  appends the signature to the key before encryption and publication.  $k_q$  is signed the same way, however, with  $\mathcal{H}$  and  $\mathcal{R}$  switched in role. Any secure signature scheme can be used, but to allow for analysis of our proposal, we assume EdDSA as the signature scheme.

Since the signature and key,  $k \parallel \sigma$ , are encrypted with the public key of the board,  $\mathcal{H}$  and  $\mathcal{R}$  are unable to inspect the context of the ciphertext after encryption. This can be solved by comparing ciphertexts. Given key pair  $\langle pk, sk \rangle$  and ciphertext  $c'$  generated by  $\mathcal{R}$  using deterministic encryption method  $\mathcal{X} : \mathcal{M} \times \mathcal{K} \rightarrow \mathcal{C}$ ,  $\mathcal{H}$  can obtain  $c$  encrypting message  $m$  using key  $pk$ . If  $c' = c$ ,  $\text{dec}_{sk}(c') = \text{dec}_{sk}(c) = m$ . However, deterministic encryption methods are not IND-CPA secure. Therefore,  $\mathcal{H}$  and  $\mathcal{R}$  use  $s$  as a seed to generate cryptographically secure random numbers used for masking the encryption. Any IND-CCA secure hybrid encryption scheme can be used. We use ECIES in our proposal for the complexity analysis.

After the first authenticity test, the keys are encrypted and published. The structure of the key communication blocks is shown in Table 4.3.

**Step 3: Verification of the keys.** The ciphertexts for both keys are retrieved by  $\mathcal{H}$  and  $\mathcal{R}$  from the distributed ledger. When both the signature of the keys and the ciphertext of the publication are valid, the keys can be used for communication. If either the signature or the ciphertext was incorrect, a deviation from the protocol has been detected and the key generation is aborted. At abortion, the **ABORT**block is publicized. Communication should continue only if both parties are certain that the Ledger can recover the correct communication keys. In the encryption of the keys, a supplementing technique is required. Deterministic encryption allows the other party to trivially verify the authenticity of the encrypted communication key.

**Aborting.** After generation of the communication keys both  $\mathcal{H}$  and  $\mathcal{R}$  verify the validity of the counterparts block. The validation assures that the other party has indeed published the correct key. If the validation succeeds, the parties can start sending queries and results. Otherwise, a special **ABORT**block is published, invalidating the generated keys.

**Querying.** A research group can communicate a query to the hospital by publicizing a Query Block on the ledger. The Query Block includes a reference to the **QUERY\_KEY** Block, the identifiers of  $\mathcal{H}$  and  $\mathcal{R}$  and the encrypted query. The structure of the Query Block is shown in Table 4.4.

**Responding.** A hospital that receives a Query Block decrypts the query

and evaluates this on its dataset. The primary keys of the records that are being used for the query are stored by the hospital and additionally aggregated into a single hash. Differential privacy is applied to the result to hide the inclusion of individual patients. This adapted result is publicized on the ledger together with the hash of the primary keys. See Algorithm 3 for the evaluation and Table 4.5 for the structure of the Response Block.

**Key recovery.** The board can recover either  $k_q$  or  $k_r$  using the encrypted forms of the keys stored on the ledger.

---

**Algorithm 2** The generation of communication keys used by  $\mathcal{H}$  and  $\mathcal{R}$ .
 

---

| $\mathcal{H}$  |                               | $\mathcal{R}$  |
|--|-------------------------------|--|
| <b>Input</b><br>$id_{\mathcal{H}}, id_{\mathcal{R}}, pk_{\mathcal{R}}, pk_{\mathcal{H}} = g^{sk_{\mathcal{H}}}$  | $g$                           | <b>Input</b><br>$id_{\mathcal{H}}, id_{\mathcal{R}}, pk_{\mathcal{H}}, pk_{\mathcal{R}} = g^{sk_{\mathcal{R}}}$  |
| <b>Output</b><br>$k_r, k_q, \sigma_q$  |                               | <b>Output</b><br>$k_r, k_q, \sigma_r$  |
| $x_q \xleftarrow{R} \mathbb{Z}_q, x_r \xleftarrow{R} \mathbb{Z}_q, x_s \xleftarrow{R} \mathbb{Z}_q$<br>$X_q \leftarrow g^{x_q}, X_r \leftarrow g^{x_r}, X_s \leftarrow g^{x_s}$  |                               | $y_q \xleftarrow{R} \mathbb{Z}_q, y_r \xleftarrow{R} \mathbb{Z}_q, y_s \xleftarrow{R} \mathbb{Z}_q$  |
|  | $\xrightarrow{X_q, X_r, X_s}$ |  |
|  | $\xleftarrow{Y_q, Y_r, Y_s}$  |  |
| $k_q \leftarrow \text{hkdf}(Y_q^{x_q}, (pk_{\mathcal{R}})^{x_q}, Y_q^{sk_{\mathcal{H}}}, X_q, Y_q, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  |                               | $Y_q \leftarrow g^{y_q}, Y_r \leftarrow g^{y_r}, Y_s \leftarrow g^{y_s}$<br>$k_q \leftarrow \text{hkdf}(X_q^{y_q}, X_q^{sk_{\mathcal{R}}}, (pk_{\mathcal{H}})^{y_q}, X_q, Y_q, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  |
| $k_r \leftarrow \text{hkdf}(Y_r^{x_r}, (pk_{\mathcal{R}})^{x_r}, Y_r^{sk_{\mathcal{H}}}, X_r, Y_r, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  |                               | $k_r \leftarrow \text{hkdf}(X_r^{y_r}, X_r^{sk_{\mathcal{R}}}, (pk_{\mathcal{H}})^{y_r}, X_r, Y_r, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  |
| $s \leftarrow \text{rng}(Y_s^{x_s}, (pk_{\mathcal{R}})^{x_s}, Y_s^{sk_{\mathcal{H}}}, X_s, Y_s, pk_{\mathcal{H}}, pk_{\mathcal{R}})$   |                               | $s \leftarrow \text{rng}(X_s^{y_s}, X_s^{sk_{\mathcal{R}}}, (pk_{\mathcal{H}})^{y_s}, X_s, Y_s, pk_{\mathcal{H}}, pk_{\mathcal{R}})$   |
| $\sigma_q \leftarrow \text{sign}_{sk_{\mathcal{H}}}(k_q)$  | $\xrightarrow{\sigma_q}$      | $\sigma_r \leftarrow \text{sign}_{sk_{\mathcal{R}}}(k_r)$  |
|  | $\xleftarrow{\sigma_r}$       |  |
| <b>if</b> verify( $pk_{\mathcal{R}}, c_r, \sigma_r$ ) $\neq 1$ <b>then</b><br><b>return</b> $\langle \perp, \perp \rangle$<br><b>end if</b>  |                               | <b>if</b> verify( $pk_{\mathcal{H}}, c_q, \sigma_q$ ) $\neq 1$ <b>then</b><br><b>return</b> $\langle \perp, \perp \rangle$<br><b>end if</b>  |
| $r_r \leftarrow s, r_s \leftarrow s$<br>$c_r \leftarrow \text{PK. enc}_{pk_{\mathcal{B}}}(k_r \parallel \sigma_r, r_r)$<br>$c_q \leftarrow \text{PK. enc}_{pk_{\mathcal{B}}}(k_q \parallel \sigma_q, r_q)$<br>pub (RESULT_KEY, $id_h, id_r, c_r$ )<br>$p, c'_q \leftarrow \text{retrieve}(\text{QUERY\_KEY}, id_h, id_r)$<br><b>if</b> $c'_q \neq c_q$ <b>then</b><br>pub (ABORT, $id_h, p$ )<br><b>return</b> $\langle \perp, \perp \rangle$<br><b>end if</b><br><b>return</b> $\langle k_r, k_q \rangle$ |                               | $r_r \leftarrow s, r_s \leftarrow s$<br>$c_q \leftarrow \text{PK. enc}_{pk_{\mathcal{B}}}(k_q \parallel \sigma_q, r_q)$<br>$c_r \leftarrow \text{PK. enc}_{pk_{\mathcal{B}}}(k_r \parallel \sigma_r, r_r)$<br>pub (QUERY_KEY, $id_h, id_r, c_q$ )<br>$p, c'_r \leftarrow \text{retrieve}(\text{RESULT\_KEY}, id_h, id_r)$<br><b>if</b> $c'_r \neq c_r$ <b>then</b><br>pub (ABORT, $id_r, p$ )<br><b>return</b> $\langle \perp, \perp \rangle$<br><b>end if</b><br><b>return</b> $\langle k_r, k_q \rangle$ |

---

---

**Algorithm 3** Evaluating the query.

---

**Input** $c$  encrypted query $k_q$  query key $k_r$  response key**Output** $r'$  response $h$  hash of all primary keys $q \leftarrow \text{dec}_{k_q}(c)$  $\langle pks, r \rangle \leftarrow \text{eval}(q)$  $r' \leftarrow \text{dp}(r, \Delta f)$  $h \leftarrow \text{hash}(pks)$ **return**  $\langle r', h \rangle$ 


---

## Chapter 5

# Evaluation

In this chapter, we will evaluate the proposed protocol in terms of complexity and security. Additionally, we will analyze the performance of a simulated run. For the complexity and security analysis, we focus on the generation of communication keys and the recovery of the communication keys.

### 5.1 Complexity Analysis

To analyze the performance of our protocol, we analyze two type of complexities: Computational complexity and message complexity. We will evaluate our protocol in three dimensions: Computational complexity, and message complexity.

#### 5.1.1 Computational Complexity

The computational complexity denotes the number of elementary operations required to run the protocol. We use group exponentiation operations as a metric for the time complexity of the algorithm. An overview of the computational complexity analysis is shown in Table 5.1.

Table 5.1: Computational complexity analysis of the protocol. Expressed in exponentiation operations in  $G_1$ , and  $G_2$ .

| Step                         | Complexity                       |
|------------------------------|----------------------------------|
| Setup party                  | $\mathbf{G}_1 + \mathbf{G}_2$    |
| Communication Key Generation | $14\mathbf{G}_1 + 4\mathbf{G}_2$ |
| Key recovery                 | $2\mathbf{G}_1 + 2\mathbf{G}_2$  |
| Request/Response             | 0                                |

**Communication Key Generation.** We analyze the computational complexity of the communication key generation shown in Algorithm 2, using

the proposed cryptographic schemes for signing and public key encryption. The computational complexity is for a single party, both  $\mathcal{H}$  and  $\mathcal{R}$  have the same computational cost for generating the communication keys. The computational complexity of the individual components is shown in Table 5.2.

Table 5.2: The different building blocks used in Algorithm 2 and their computational complexity. Expressed in exponentiation operations in  $G_1$ , and  $G_2$ .

| Algorithm     | Times performed in Communication Key Generation | Complexity      |
|---------------|---|-----------------|
| Triple DH     | 3   | $4\mathbf{G}_1$ |
| EdDSA sign    | 1   | $2\mathbf{G}_2$ |
| EdDSA verify  | 1   | $2\mathbf{G}_2$ |
| ECIES encrypt | 2   | $2\mathbf{G}_1$ |

We can see from Table 5.1 that all computational complexities are constant time. The number of board members has no influence on the complexity of the algorithms. This makes the protocol very scalable; it can be employed with any number of parties without impacting the performance.

### 5.1.2 Message Complexity

The protocol is reliant on communication over authenticated channels and a distributed ledger. These channels are equivalent in cost, and we will therefore separate them in the message complexity analysis. We use  $L$  to denote messages sent over the distributed ledger and  $D$  for messages sent over an authenticated channel. The message complexity is shown in Table 5.3.

Note that the key recovery will require  $t^2$  messages over authenticated channels, where  $t$  is the threshold of DKG. Since DKG requires  $t$  parties to recover the secret, we need  $t^2$  for every party to share each partial decryption with one another.

Table 5.3: Message complexity analysis of the protocol.  $L$  denotes messages sent over the distributed ledger and  $D$  for messages sent over an authenticated channel.

| Step                         | Complexity |
|------------------------------|------------|
| Setup party                  | $L$        |
| Communication Key Generation | $4D + 4L$  |
| Key recovery                 | $t^2D$     |
| Request/Response             | $L$        |



## 5.2 Security Analysis

### 5.2.1 Correctness of the Communication Keys Generation

We prove that after correct execution of 2,  $\mathcal{R}$  and  $\mathcal{H}$  will end up with the same communication keys.

We take first the generation of the query key. Let  $k_q$  be the query key generated by  $\mathcal{H}$  and  $k'_q$  the query key generated by  $\mathcal{R}$ . The query key is correctly generated if after execution of the algorithm  $k_q = k'_q$ .  $\mathcal{H}$  and  $\mathcal{R}$  calculate  $k_q \leftarrow \text{hkdf}(Y_q^{x_q}, (pk_{\mathcal{R}})^{x_q}, Y_q^{sk_{\mathcal{H}}}, X_q, Y_q, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  and  $k'_q \leftarrow \text{hkdf}(X_q^{y_q}, X_q^{sk_{\mathcal{R}}}, (pk_{\mathcal{H}})^{y_q}, X_q, Y_q, pk_{\mathcal{H}}, pk_{\mathcal{R}})$  respectively. Since  $\text{hkdf}$  is a deterministic function,  $k_q$  and  $k'_q$  are equal if all arguments are the same;

$$Y_q^{x_q} = X_q^{y_q}, \quad (5.1)$$

$$(pk_{\mathcal{R}})^{x_q} = X_q^{sk_{\mathcal{R}}}, \quad (5.2)$$

$$Y_q^{sk_{\mathcal{H}}} = (pk_{\mathcal{H}})^{y_q}, \quad (5.3)$$

$$X_q = X_q, \quad (5.4)$$

$$Y_q = Y_q, \quad (5.5)$$

$$pk_{\mathcal{H}} = pk_{\mathcal{H}}, \quad (5.6)$$

$$pk_{\mathcal{R}} = pk_{\mathcal{R}}. \quad (5.7)$$

We assume  $pk_{\mathcal{H}}$  and  $pk_{\mathcal{R}}$  to be shared correctly prior to the execution of Algorithm 2. Additionally, since we assume an authenticated channel,  $X_q$  and  $Y_q$  will be equal for both parties. This means that Eq. (5.4) to Eq. (5.7) hold.

Eq. (5.1) holds since

$$Y_q^{x_q} = (g^{y_q})_q^x = g^{y_q \cdot x_q} = (g^{x_q})_q^y = X_q^{y_q}. \quad (5.8)$$

Similarly Eq. (5.2) and Eq. (5.3) hold as well;

$$(pk_{\mathcal{R}})^{x_q} = g^{sk_{\mathcal{R}} \cdot x_q} = X_q^{sk_{\mathcal{R}}}, \quad (5.9)$$

$$(pk_{\mathcal{H}})^{y_q} = g^{sk_{\mathcal{H}} \cdot y_q} = Y_q^{sk_{\mathcal{H}}}. \quad (5.10)$$

Therefore, after execution  $k_q = k'_q$  and both  $\mathcal{H}$  and  $\mathcal{R}$  are in possession of a valid query key. A similar proof can be used for the response key  $k_r$ .

### 5.2.2 Proof of Auditability

We will now prove that auditability will always be achieved, under the condition that at least one of the involved members acts honest. In other words: Either  $\mathcal{H}$  and  $\mathcal{R}$  collude together, or the all query-response pairs are auditable.

Both parties publish a single communication key;  $\mathcal{H}$  publishes  $PK.\text{enc}(k_r \parallel \sigma_r)$  and  $\mathcal{R}$  publishes  $PK.\text{enc}(k_q \parallel \sigma_q)$ . The encryption is a non-deterministic, CCA secure algorithm. However,  $\mathcal{H}$  and  $\mathcal{R}$  share a seed for creating the nonce. Under a known nonce, the encryption becomes deterministic. Therefore, both parties are able to retrieve the same ciphertexts. This puts a security mechanism in place on the authenticity of the key.

We will prove this mechanism to be secure. This one party to be malicious, this protection can only be circumvented if the adversary finds a collision for the ciphertext  $c_r$ ;  $c'_r$  where Given  $m_1, m_2$  where  $m_1 \neq m_2$ , the adversary must find

$$\text{enc}(k_1, m_1) = \text{enc}(k_2, m_2) \quad (5.11)$$

for some  $k_1, k_2 \in \mathbf{K}$ .

We claim that there exist no collision for the ciphertext  $c_r$ . Given the correctness requirement of a computationally secure algorithm,

$$\text{dec}(k, \text{enc}(k, m)) = m, \quad (5.12)$$

every ciphertext has a unique decryption. Thus, there exist no ciphertext which can be decrypted into two different messages.

If at least one of the two parties is non-malicious, then publication of an incorrect ciphertext will result in the communication key generation aborted. Also, requests or responses which are indecipherable by the key agreed upon will trigger an abort. Preventing a party from publishing the abort will not compromise auditability since the party will not publish any new messages before the abort.

Combining the verifiable key publication, verifiable decryption of queries and responses, and the non-collision of ciphertexts, the auditability is guaranteed.

### 5.2.3 Proof of Key Authenticity

We prove that the board that the board is guaranteed authentic keys.

Assume  $\mathcal{H}$  tries to publish a different key from the shared key. Governing board members have access to the `RESULT_KEY` block published by  $\mathcal{H}$ . A consortium of the governing board can decrypt this ciphertext to recover  $k'_r$  and  $\sigma'_r$ . If  $\text{hash}(k'_r \parallel \sigma'_r)$  is unequal to  $\text{hash}(k_r \parallel \sigma_r)$ , then  $\mathcal{H}$  had published a malformed query key.

If the hash is correct, the consortium knows with high probability that  $\mathcal{H}$  and  $\mathcal{R}$  have agreed on  $k_q$ . For this, we give the following argumentation.

In order to publish any  $k_m$ , where  $k_m \neq k_q$ ,  $\mathcal{H}$  must construct a  $\sigma_m$  such that  $\sigma_m$  is a valid signature for  $k_m$ . However, given that sign is a strongly secure signature scheme, constructing a valid signature is only feasible with knowledge of  $sk_{\mathcal{R}}$ . Therefore, with high probability,  $k'_q$  is the shared key  $k_q$ .

When  $\mathcal{H}$  and  $\mathcal{R}$  collude, they can replace  $k_q$  or  $k_r$  by any value of their choice, without making this detectable in the `QUERY_KEY` or `RESULT_KEY` block.

Combining the fact both members act honestly, the `RESULT_KEY` and `QUERY_KEY` block will be correct.

#### 5.2.4 Security Sketch of Communication Keys

The protocol needs to guarantee the security of queries and responses. Since these are encrypted using IND-CCA secure encryption schemes, the security is dependent on the secrecy of the communication keys. We give a security sketch that suggests the communication key generation is indeed secret.

We can analyze the advantage of a passive observer in the following way. The adversary observes two separate processes. Since the generation of the query and the generation of the result key are symmetrical, we will focus on the security of a single key exchange first and combine this with possible interpolation methods to analyze the information gain of the generation of two keys.

The first observed process is the exchange of ephemeral keys between  $\mathcal{H}$  and  $\mathcal{R}$ . This provides an eavesdropper with  $X_q, Y_q, X_r, Y_r, X_s$  and  $Y_s$ . The security of this exchange is based on computational Diffie-Hellman assumption[31].

The second observed process is the sharing of the signatures. The observer receives the signature of  $k_q$  and  $k_r$ ,  $\sigma_q$  and  $\sigma_r$  respectively. The key here is protected by the pre-image resistance property of the hashing algorithm used for signing.

Since the signature scheme used is a strongly binding secure signature scheme, it is infeasible to gain information about the underlying key through the signature in polynomial time. However, it does provide the observer with an oracle to verify a correct key.

If normal execution of the key generation is followed, and the published keys are not rejected by an `ABORT`, the observer obtains  $c_q = \text{PK}.\text{enc}_{pk_B}(k_q \parallel \sigma_q, r_q)$ . Retrieving the key from  $c_q$  involves breaking the one-way property of the public key encryption scheme used. Since we use a computationally secure public key encryption scheme, the scheme is secure against message recovery by a polynomial time-bounded adversary.

Since a share of the board key does not provide any information on the actual board  $sk$ , we can reduce the view of the board members to the view of an eavesdropper. Therefore, no polynomial bounded adversary outside the involved  $\mathcal{H}$  and  $\mathcal{R}$  can learn the secret communication key.

## Chapter 6

# Discussion

Enabling the auditing of private data communication is a first step in a process which is transparent and causes trust. In medical data sharing, a crucial step in moving medical science forward, the system ensures secure communication between data owners and data processors. But more importantly, it ensures secure communication while making auditing possible, demonstrating the willingness of the parties involved to indulge in an ethical and legal approach to data analysis.

Compared to related work, our approach is unique. PPDC and VeriDC are similarly suited for medical data analysis, but provide privacy-preserving proofs to auditors. Although the data is protected, this significantly increases computational complexity since it relies on NIZK proofs. zKledger also relies on zero-knowledge proofs and requires the data and auditing to be of numerical nature. The security of MeDShare’s auditing depends on the logging of traces in the system, as opposed to zero knowledge auditing or ensuring trust through a multi-eyed approach. Additionally, trust guarantees on the auditor are not discussed.

Of course, use of the system does not guarantee that data is solely shared using the system. The system cannot prevent two parties from sharing data in outside the protocol, while one of the two parties has full ownership over the data. However, we argue that this is unavoidable.

If a data owner has direct control over data, without other assumptions of systems in place, it can delete it, copy it, alter it without any traces. The only way to know for sure that the data is authentic is to have an independent source confirming the fact. Ensuring the existence of such a source is difficult, if not outright infeasible, in many situations. The independent source should confirm every mutation on the data. From the creation, to the aggregation to finally the reading of the data. Practically, this means that a test performed by a medical caregiver in a hospital should be overseen by an independent party. Such measures are radical, expensive, and are unlikely to bring any benefit. However, for automated process such as data collec-

tion from IoT devices, this will more likely be possible. Data can be signed before being sent to the hospital, and more security practices can be made in place to ensure the authenticity of the data. This gives the possibility in transforming the role of the hospital from a data owner into a sort of data mediator, enabling patients to forward their information through the hospital to a medical health data space. Still, for now to view the hospital as the data owner reflects the current situation better.

With direct control of the data and access to the internet, it becomes possible to share the data. A party with access to private medical data could share it with other parties without the right legal precautions. We argue that the proposed system aids in prevention and detection of this. First, if a standard for data sharing is in place, adherence to this standard shows goodwill. Even if violations of data protection acts are made on the ledgers, the party still committed to publishing them on the ledger. Second, if the protocol becomes standard, deviations from the protocol become rarer, and thus easier to identify. A large data stream from a hospital could be detected by an ISP and, if not correlated with an entry in the ledger, flagged as suspicious and potential for further investigation. Additionally, if research groups cannot disclose the origin of the data on the ledger, it might be a sign for further investigation. For example, if the paper published which uses private medical data from a source, but the ledger cannot explain the data, it could be a sign of fraud.

The protocol described in this thesis has been inspired by SEPTON and TEHDAS. These projects are the manifestations of the wish for future-proof medical healthcare by the European Union. This protocol makes assumptions about the existence of significant resources and systems, however, these are feasible in a European wide context.

The European Union has the capacity and technology for the deployment of a European wide key infrastructure for hospitals and research groups. The EU requires members to assign Data Protection Authorities to guarantee national compliance with the respective implementation of the GDPR. These are good examples of parties which could function as board members in our system. Additionally, non-profit or for-profit organizations, representative groups from the healthcare, or overseers from law or medical faculties can be suitable choices for board members. The system would exploit the knowledge and skills of different parties, while guaranteeing security of the data transfers.

In our proposal, we give a simple example of combining auditable data transfers with differential privacy to limit privacy loss of the patients when data is transferred to a research group. This also prevents the hospital from losing control over the patient's data; The research group has no access to the original data. Differential Privacy, however, can be broken when the same query is performed to retrieve multiple samples. By comparing the results, more information about the value of the query can be learned. This

could be achieved, for example, by multiple research groups sending the same query and later collaborating on recovering more information about the true query result. Differential privacy might be sufficient as a privacy protection mechanism in some scenarios, however, this is insufficient in non-numerical use cases.

Thus far in this thesis, we have made the assumption that hospitals have the capacity to perform queries by research groups on demand. In general, hospitals are assumed to not have direct access to significant computing power, and work on medical data analytics often optimize or offload the hospital's computations [17][18][32]. These solutions protect sensitive medical information through privacy-preserving analytics, however, they do not provide auditing. A deployed health research space, like the proposed EHDS, should pay heed to the concerns put forth by the citizen surveys by TEHDAS and the medical professionals, and implement privacy-preserving techniques, auditing and minimize the computational complexity for hospitals.

## Chapter 7

# Future Work

In this chapter, we discuss some possible future directions for auditable and secure data sharing. For real world deployment of auditable and secure data sharing, combining secure computation with delegating analysis is required. Additionally, using auditable and secure data sharing for reproducible medical research or applications outside of medical research.

### 7.1 Privacy-preserving Outsourcing of Data Analytics

Integrating privacy-preserving outsourcing of data analytics can significantly reduce the computing load on hospitals. In our proposal, hospitals are required to run the queries. The work of Blanton et al. [17] demonstrated a method of privacy-preserving learning from distributed medical records through secure multi-party computations. Hospitals pre-process sensitive data and encrypt the data in a special way. The encrypted data are split and each computational service provider receives a share. The computational service provider can perform analytics on the data to get an encrypted result. These results are gathered by a research group who by combining the encrypted data can break the encryption and learn the result. By extending our protocol with an additional party type, computational service provider, the privacy-preserving learning algorithm by Blanton et al. can become auditable. When a colloquium of board members recovers all data sent to the computational service providers, the privacy-preserving mechanism can be broken and the original data used for training can be inspected. The limitation that computational service providers are not allowed to collude, can actually be turned into a positive property in regard to auditing.

## 7.2 Reproducible Medical Research

Our proposal connects a history of queries and results together. By linking publications to the relevant queries and results, there exist an opportunity for improving the reproducibility of medical research. If reviewers receive access to these transcriptions, it would aid them in reviewing the reported procedure and soundness of data. Additionally, at publication, relevant information could be disclosed to the public. For our proposal, results and queries cannot be released separately. Since communication keys are designed to be reused, research groups would have to decide in advance which queries to encrypt with the same keys. This prevents the research group from picking queries and results to be released to the public. A system more tailored towards this use case could additionally implement ways to track publications on the ledger.

## 7.3 Further Auditable Data Sharing

Our protocol focuses on making data sharing between Hospitals and Research groups auditable. The same protocol can be applied to other data sharing algorithms to make them auditable. We can abstract the party requirements for the algorithm. We can use the scenario when two parties communicate sensitive data which needs to adhere to a standard, which possibly might not be verifiable by the receiving parties. Further use cases, outside the medical domain, can be explored.



## Chapter 8

# Conclusion

In this thesis, we demonstrate a secure and efficient method for distributed auditing of data sharing. We show that auditing of secure medical data sharing can be achieved in a distributed fashion using, a decentralized ledger, distributed key generation and recoverable key generation. The proposed system guarantees auditability and secrecy under a realistic threat model. Additionally, the overhead for making the protocol auditable is limited, all algorithms are performed in constant time. Previous work has demonstrated that DKG will scale poorly with the number of involved board members. Also, the dependency on tracking all queries and responses does introduce overhead and leaves opportunity for future work to explore efficient ways of sharing data over distributed ledgers. However, when a DKG and communication over distributed ledgers is efficient enough for real world scenarios, our protocol introduces little overhead to make it auditable.

The protocol provides secrecy and auditability using the described cryptographic algorithms for secure signing and authenticated encryption. However, these cryptographic algorithms can be switched out for other primitives without loss of security. Depending on the scenario, different elliptic-curves can be chosen to increase the performance.

# Bibliography

- [1] J. Scheibner, J. L. Raisaro, J. R. Troncoso-Pastoriza, *et al.*, “Revolutionizing medical data sharing using advanced privacy-enhancing technologies: Technical, legal, and ethical synthesis,” *JMIR*, vol. 23, no. 2, e25120, 2021. DOI: 10.2196/25120.
- [2] J. C. Wong and O. Solon, “NHS cyber-attack: GPs and hospitals hit by ransomware,” *BBC NewsHealth*, May 12, 2017. [Online]. Available: <https://www.bbc.com/news/health-39899646> (visited on 05/10/2023).
- [3] Nederlandse Zorgautoriteit. “Zorgprestatie Model,” Nederlandse Zorgautoriteit. (2023), [Online]. Available: <https://zorgprestatiemodel.nza.nl/> (visited on 05/08/2023).
- [4] R. Huissen. “Vraag & Antwoord,” Vertrouwen in de GGZ. (2023), [Online]. Available: <https://vertrouwenindeggz.nl/faq> (visited on 05/08/2023).
- [5] K. Menager, J. Maddocks, L. Mathieu, R. Richards, M. Saelaert, and W. Van Hoof, “TEHDAS study to assess citizens perception of sharing health data for secondary use,” Mar. 31, 2023.
- [6] R. Wang, W. Tsai, J. He, C. Liu, Q. Li, and E. Deng, “A medical data sharing platform based on permissioned blockchains,” in *ICBTA*, ACM, 2018, pp. 12–16. DOI: 10.1145/3301403.3301406.
- [7] M. Attaran, “Blockchain technology in healthcare: Challenges and opportunities,” *Int. J. Healthc. Manag.*, vol. 15, no. 1, pp. 70–83, Jan. 2, 2022, ISSN: 2047-9700, 2047-9719. DOI: 10.1080/20479700.2020.1843887. (visited on 02/06/2023).
- [8] J. Gluck, F. Schaub, A. Friedman, *et al.*, “How short is too short? implications of length and framing on the effectiveness of privacy notices,” in *SOUPS*, USENIX Association, 2016, pp. 321–340.
- [9] A. M. McDonald and L. F. Cranor, “The cost of reading privacy policies,” *I/S: A Journal of Law and Policy for the Information Society*, vol. 4, pp. 540–565, 2008.

- [10] E. Fujisaki and T. Okamoto, “Secure integration of asymmetric and symmetric encryption schemes,” in *CRYPTO*, M. J. Wiener, Ed., ser. Lecture Notes in Computer Science, vol. 1666, Springer, 1999, pp. 537–554. DOI: 10.1007/3-540-48405-1\_34.
- [11] E. Fujisaki and T. Okamoto, “Secure integration of asymmetric and symmetric encryption schemes,” *J. Cryptol.*, vol. 26, no. 1, pp. 80–101, 2013. DOI: 10.1007/s00145-011-9114-1.
- [12] T. Okamoto, E. Fujisaki, and H. Morita, “Psec: Provably secure elliptic curve encryption scheme (submission to p1363a),” IEEE P1363a, Tech. Rep., 1999.
- [13] A. Menezes, “Evaluation of security level of cryptography: The revised version of psec-2 (psec-kem),” Technical report, CRYPTREC, Tech. Rep., 2001.
- [14] A. W. Dent, “Ecies-kem vs. psec-kem,” Technical Report NES/DOC/RHU/WP5/028/2, NESSIE, Tech. Rep., 2002.
- [15] V. G. Martínez, L. H. Encinas, *et al.*, “A comparison of the standardized versions of ecies,” in *2010 Sixth International Conference on Information Assurance and Security*, IEEE, 2010, pp. 1–4. DOI: 10.1109/isias.2010.5604194.
- [16] T. P. Pedersen, “A threshold cryptosystem without a trusted party (extended abstract),” in *EUROCRYPT*, D. W. Davies, Ed., ser. Lecture Notes in Computer Science, vol. 547, Springer, 1991, pp. 522–526. DOI: 10.1007/3-540-46416-6\_47.
- [17] M. Blanton, A. R. Kang, S. Karan, and J. Zola, “Privacy preserving analytics on distributed medical data,” *CoRR*, vol. abs/1806.06477, 2018. arXiv: 1806.06477.
- [18] X. Liu, Y. Zheng, X. Yuan, and X. Yi, “Medisc: Towards secure and lightweight deep learning as a medical diagnostic service,” in *ESORICS*, E. Bertino, H. Shulman, and M. Waidner, Eds., ser. Lecture Notes in Computer Science, vol. 12972, Springer, 2021, pp. 519–541. DOI: 10.1007/978-3-030-88418-5\_25.
- [19] Q. Xia, E. B. Sifah, K. O. Asamoah, J. Gao, X. Du, and M. Guizani, “Medshare: Trust-less medical data sharing among cloud service providers via blockchain,” *IEEE Access*, vol. 5, pp. 14 757–14 767, 2017. DOI: 10.1109/ACCESS.2017.2730843.
- [20] D. Bonyuet, “Overview and impact of blockchain on auditing,” *International Journal of Digital Accounting Research*, vol. 20, pp. 31–43, 2020. DOI: 10.4192/1577-8517-v20\_2.

- [21] K. Fanning and D. P. Centers, “Blockchain and its coming impact on financial services,” *Journal of Corporate Accounting & Finance*, vol. 27, no. 5, pp. 53–57, 2016. DOI: 10.1002/jcaf.22179. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jcaf.22179>.
- [22] J. L. Alarcon and C. Ng, “Blockchain and the future of accounting,” *Pennsylvania CPA Journal*, vol. 88, no. 4, pp. 26–29, 2018.
- [23] Y. Bakos, H. Halaburda, and C. Müller-Bloch, “When permissioned blockchains deliver more decentralization than permissionless,” *Commun. ACM*, vol. 64, no. 2, pp. 20–22, 2021. DOI: 10.1145/3442371.
- [24] Y. Chen, S. Ding, Z. Xu, H. Zheng, and S. Yang, “Blockchain-based medical records secure storage and medical service framework,” *Journal of Medical Systems*, vol. 43, pp. 1–9, 2019. DOI: 10.1007/s10916-018-1121-4.
- [25] N. Narula, W. Vasquez, and M. Virza, “Zkledger: Privacy-preserving auditing for distributed ledgers,” in *NSDI*, S. Banerjee and S. Seshan, Eds., USENIX Association, 2018, pp. 65–80.
- [26] Y. Zhao, X. Yang, Y. Yu, B. Qin, X. Du, and M. Guizani, “Blockchain-based auditable privacy-preserving data classification for internet of things,” *IEEE Internet Things J.*, vol. 9, no. 4, pp. 2468–2484, 2022. DOI: 10.1109/JIOT.2021.3097890.
- [27] X. Zheng, Y. Zhao, H. Li, R. Chen, and D. Zheng, “Blockchain-based verifiable privacy-preserving data classification protocol for medical data,” *Comput. Stand. Interfaces*, vol. 82, p. 103605, 2022. DOI: 10.1016/j.csi.2021.103605.
- [28] S. Xu, X. Cai, Y. Zhao, *et al.*, “Zkrpchain: Towards multi-party privacy-preserving data auditing for consortium blockchains based on zero-knowledge range proofs,” *Future Gener. Comput. Syst.*, vol. 128, pp. 490–504, 2022. DOI: 10.1016/j.future.2021.09.034.
- [29] A. Kate, Y. Huang, and I. Goldberg, “Distributed key generation in the wild,” *IACR*, p. 377, 2012. DOI: 10.1109/icdcs.2009.21.
- [30] M. Marlinspike and T. Perrin, “The x3dh key agreement protocol,” *Open Whisper Systems*, vol. 283, p. 10, 2016.
- [31] D. Boneh and V. Shoup, “A graduate course in applied cryptography,” *Draft 0.6*, 2023.
- [32] C. Bonte, E. Makri, A. Ardeshirdavani, J. Simm, Y. Moreau, and F. Vercauteren, “Towards practical privacy-preserving genome-wide association study,” *BMC Bioinform.*, vol. 19, no. 1, 537:1–537:12, 2018. DOI: 10.1186/s12859-018-2541-3.