



Crypto Bluebook

Version 1.3.2

Contents

I	Personal Health Data Platform	4
1	Introduction	5
1.1	End-to-end encryption	5
1.2	Document structure	5
2	Cryptographic basics	7
2.1	Symmetric encryption	7
2.2	Asymmetric encryption	11
2.3	Hybrid encryption	12
3	Secure keys	14
3.1	Keys overview	14
3.2	Account creation	17
3.3	Login	19
3.4	Account recovery	20
3.5	Key maps	21
4	Data model and data access	23
5	Data ingestion	27
5.1	Onboarding	27
5.2	Document upload	28
5.3	Document access by client	29
6	Data donation	30
6.1	Requirements	30
6.2	Design Decisions	31
6.3	Register as a donor	31
6.4	Donate a document	35
6.5	Revocation	37
II	Formal treatment	41
7	Cryptographic notation	42

8 Account creation	46
III Appendices	49
A Bibliography	50
B Glossary	52
C Change history	54
D Acknowledgments	55

List of figures

2.1 Symmetric encryption and decryption	8
2.2 Asymmetric encryption and decryption	12
2.3 Hybrid encryption and decryption	13
3.1 Key encryption relationships.	15
3.2 Keys and key ciphertexts involved to decrypt a document	16
3.3 Download of recovery key during account creation	18
3.4 Example of key maps	22
4.1 Data model example.	24
4.2 Data layout for a encrypted records and attachments	25
4.3 Data decryption flow for reading a record	26
4.4 Data decryption flow for reading an attachment	26
6.1 User registers as donor.	33
6.2 User donates document	36
6.3 User revokes donation consent.	38
8.1 User account creation data	48

List of tables

3.1	Encryption relationships of different key types	16
3.2	Data elements of account creation	19
4.1	Data elements stored for a record	24
5.1	Ciphertexts sent to the server by two parties	28
7.1	Different types of cryptographic symbols.	43
7.2	Superscripts indicating subtypes of cryptographic symbols.	43
7.3	Examples of the cryptographic notation.	45
8.1	User data required for account creation	46
8.2	Automatically generated data by the client.	47
8.3	User registration payload that gets sent to the server.	47

Part I

Personal Health Data Platform

1 Introduction

Data4Life develops and operates a data platform called Personal Health Data Platform (PHDP) that uses end-to-end encryption (E2EE) to allow users to securely store and access healthcare data, receive healthcare data from external sources (like hospitals), and selectively share data with third parties (like doctors or other healthcare professionals). This document describes the cryptographic protocols¹ that implement the above-mentioned data-related tasks.

1.1 End-to-end encryption

Before users can access our platform, they must first register and validate an account with us. This establishes a private and end-to-end encrypted data storage that can be accessed by the applications we offer. By an *application* we refer to either a browser-based application or a native mobile app. In most cases we will collectively refer to them as *client applications*, or just *clients*. We will often refer to the private data storage as just the *server*.

Users may, for example, select documents from hard disk and upload them into their data storage. The documents can later be downloaded again. At no point in time does Data4Life get access to the unencrypted data. Data is always encrypted on the client before it is transmitted to the server. Similarly, when downloading documents, the encrypted data is sent from the server to the client where it is decrypted and displayed.

In addition to manually uploading documents, users can also grant one or more third parties append-only access to their data storage. This way, a hospital may transmit discharge letters or similar documents directly into the users' data storage. Apart from granting append-only access once, no further user interaction is necessary.²

1.2 Document structure

Chapter 2 introduces cryptographic concepts required to follow the subsequent chapters:

¹We adopt the notion of a cryptographic protocol from [10] and consider them to “*consist of an exchange of messages between participants.*”

²And, of course, users can revoke that write-only access at any time.

- Chapters 3 and 4 describe the internal data model and the protocols for user data upload and user data download.
- Chapter 5 explains how third parties can securely write into users' data storage.
- Chapter 6 explains how users can donate data to the Analytics Platform (ALP).

2 Cryptographic basics

We assume the reader has enjoyed some elementary exposure to the topics of cryptography. This document is no attempt to provide a thorough treatment of the theoretical and practical underpinnings of the vast field that is cryptography. We will briefly revisit the topics that are required to follow the next chapters. For further reading we refer to the standard literature [2, 10, 17, 14, 20].

2.1 Symmetric encryption

Encryption is the principal goal of cryptography; it makes data incomprehensible in order to ensure its confidentiality [2]. The data or documents to be encrypted are also often referred to as *messages* or *plaintexts*. An algorithm called a *cipher* gets as input the plaintext and a secret (the *key*) and produces the encrypted output called the *ciphertext*.

If the same key is used for encryption and decryption, we are dealing with *symmetric encryption* and a *symmetric cipher*. Figure 2.1 illustrates this scenario. As it is customary in cryptographic literature, we use special given names when denoting the participating parties in a cryptographic protocol. As usual, Alice wants to send a secret message to Bob across an insecure channel which might be eavesdropped on by some malicious party Eve.

Both Alice and Bob share the secret key k . Alice uses the encryption function Enc of some symmetric cipher to encrypt message m with key k to produce the ciphertext C which gets sent across the insecure channel to Bob who uses the same key k and the decryption function Dec of the symmetric cipher to reproduce the message m from the ciphertext C .

Data4Life uses the Advanced Encryption Standard cipher with 256-bit keys (AES-256) for all symmetric encryption.

2.1.1 Advanced Encryption Standard (AES)

Advanced Encryption Standard (AES) is a symmetric block cipher which became a standard in 2001 [8] after it emerged as the lead candidate in a U.S. government call-for-algorithms to replace the outdated DES cipher [16]. The block size of AES is 128 bits, that is, the encryption and decryption functions of it operate on 128-bit

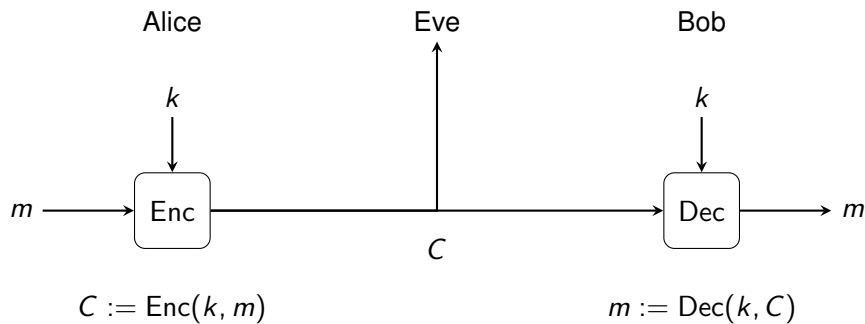


Figure 2.1: Symmetric encryption and decryption between Alice and Bob.

messages and 128-bit ciphertexts, respectively. AES supports three key lengths: 128, 192 and 256 bits. At Data4Life we exclusively use 256-bit keys and refer to it as AES-256. Hence, the Enc function (and also Dec) of Figure 2.1 has the following signature for AES-256:

$$\text{Enc}_{\text{AES-256}} : \{0, 1\}^{256} \times \{0, 1\}^{128} \rightarrow \{0, 1\}^{128}$$

Block ciphers such as AES encrypt only a fixed number of bits at a time. For longer messages (and in our case a message will be a healthcare document, possibly as large as an X-ray image) we need to iteratively apply the cipher on consecutive chunks of the message in order to fully encrypt it (and likewise for decryption). There are various ways to iteratively apply a block cipher, which is known as the *modes of operation*. Modes of operations are encryption algorithms of their own, but they need a block cipher as a “plugin” to carry out the actual encryption. Note also, that a mode of operation is not necessarily tied to a specific block cipher. We will discuss some of them below in the context of AES-256, of course, but any other block cipher would do.

2.1.2 Electronic Code Book Mode (ECB)

A straightforward, but problematic approach to extend AES, or any block cipher, to arbitrary message sizes is to split up the original message m into n chunks of size 128 bits each¹ and encrypt those chunks m_i individually:

$$\begin{aligned}
 C_i &:= \text{Enc}(k, m_i), \quad \text{for } i = 1, \dots, n \\
 m_i &:= \text{Dec}(k, C_i), \quad \text{for } i = 1, \dots, n
 \end{aligned}$$

While this idea has some desirable properties (encryption and decryption can be parallelized trivially, and random access to the ciphertext is also possible),² it maps

¹We will discuss further down how to handle the case when the message length is not an integer multiple of the block size, that is, when the last block is not full.

²Assume the message m is an ultrasound video and one wants to fast-forward to, say, the middle. One does not have to decrypt the first half of the ciphertext to reach the middle, but one can directly seek to the center block and start decrypting from there.

same message blocks to same ciphertext blocks, that is:

$$m_i = m_j \Leftrightarrow C_i = C_j$$

As a result of this block independence, each repetition of a plaintext block results in repetition of the corresponding ciphertext block, presenting a cryptographic weakness by unnecessarily revealing structural information about the message [10]. This mode of operation is not used in any Data4Life application for its obvious security compromise. It is included here only to motivate the need for different modes of operation.

2.1.3 Cipher Block Chaining Mode (CBC)

The main idea to overcome the drawback of ECB mode is to make a ciphertext block C_i not only dependent on its message block m_i and the key k , but also on the previous ciphertext block C_{i-1} and thus indirectly dependent on all previous ciphertext blocks:

$$\begin{aligned} C_i &:= \text{Enc}(k, m_i \oplus C_{i-1}), \quad \text{for } i = 1, \dots, n \\ m_i &:= \text{Dec}(k, C_i) \oplus C_{i-1}, \quad \text{for } i = 1, \dots, n \end{aligned}$$

Before encryption, each message block m_i is combined with the previous ciphertext block C_{i-1} using the XOR (\oplus) operation. This effectively randomizes the message block using the previous ciphertext block. The only question remaining is what value to pick for C_0 , which is needed to produce the first ciphertext block C_1 . This value is called an *initialization vector* and there are various strategies to pick it [10].

At Data4Life we use AES-256 in CBC mode with zero initialization vector to encrypt and decrypt so-called *tags*. The zero initialization vector is used to achieve a deterministic ciphertext. See Chapter 4 for details.

2.1.4 Padding

The two modes discussed so far essentially “keep the block nature” of the underlying block cipher. That is, a message is still encrypted and decrypted block-wise using the block cipher’s encryption and decryption function, respectively. In this scenario it can happen that the last chunk m_n of a message does not have full block size. For example, a 100-byte message would—assuming AES—result in seven blocks which all have to be 128 bits in length in order to work for AES. However, the last block would only contain 32 bits of data, leaving 96 bits unused ($800 = 6 \cdot 128 + 32$).

Some *padding* has to be applied to fill up such last blocks without compromising the security of the cipher and also to enable the receiving party to detect those additional bits and remove them in order not to confuse them with message data [10]. As we will see below, there are modes of operation that do not require padding.

2.1.5 Output Feedback Mode (OFM)

Even though Output Feedback Mode is not used at Data4Life, we use it here to motivate certain properties of GCM (which we describe in the next section), which is used to encrypt all data (except tags). OFM does not use the block cipher to directly encrypt the message. Instead it uses the encryption function Enc of the block cipher to produce a pseudorandom key stream k_i . The actual encryption is implemented by XORing the message bits (interpreted as a message stream) with the corresponding bits of the key stream. Decryption is carried out by XORing the ciphertext bits with the same key stream bits.

$$\begin{aligned} C_i &:= m_i \oplus k_i & k_i &:= \text{Enc}(k, k_{i-1}), \quad \text{for } i = 1, \dots, n \\ m_i &:= C_i \oplus k_i & k_0 &:= IV \end{aligned}$$

Since the XOR operation is carried out on bit level, no padding is required. This essentially constructs a *stream cipher* out of a block cipher. Also note that the decryption function of the underlying block cipher is never used. The receiving party must just recreate the exact same key stream to decrypt the message. Hence, the same mechanism for creating the key stream must be carried out on both sides.

The initialization vector IV must be unique for each message and must never be repeated because it determines the key stream (together with the key k , of course). The initialization vector is typically prepended to the ciphertext stream. It must be unique, but it is not kept secret.

2.1.6 Galois/Counter Mode (GCM)

Using Galois/Counter Mode [7] constructs a stream cipher from the underlying block cipher. It also has another feature which the modes discussed above did not exhibit: *authenticated encryption*. All encryption ciphers and modes discussed so far catered for one thing only: message confidentiality. However, there is no means to detect errors or malicious modifications in the ciphertexts. That is, if a bit flipped in the ciphertext, the receiver Bob would still get a result from the decryption. Only this result would be different from the original plaintext that Alice sent. However, it would be up to Bob to detect whether the decrypted message has changed from its original. From the cipher and mode algorithm point of view, everything worked out fine. After all, the XOR operation and the internal workings of AES operate on bit level and do not care about any higher-level format.

For the receiver of a ciphertext it is desirable to learn from the decryption algorithm whether every bit has made it through the communication channel unaltered. This is what authenticated encryption caters to. The output of GCM will not only contain the ciphertext, but also a so-called *authentication tag* which acts as a cryptographic checksum that can be used to detect modifications. To be more precise, GCM offers *authenticated encryption with associated data (AEAD)*. Alice, in addition to the secret message m , may also add another piece of information A , which does not get encrypted, but is covered by the authentication tag.

The intricate workings of GCM go beyond the scope of this document, but the following pseudo-code shows the “GCM API” from a programmer’s point of view:³

Encryption:	Decryption:
Input: - message m	Input: - ciphertext C
- key k	- key k
- additional data A	- additional data A
- initialization vector IV	- initialization vector IV
	- authentication tag T
Output: - ciphertext C	Output: - message m , or error
- additional data A	
- authentication tag T	
(covering IV , C and A)	

If Alice wants to send a message m and additional data A to Bob, a random initialization vector IV is chosen and the concatenation (\parallel) of the following data gets transmitted to Bob. (The additional data A , if specified, would be encoded into C and is not explicitly mentioned.)

$$IV \parallel C \parallel T$$

Again, A is protected by the authentication tag, but not encrypted. At Data4Life we do not use the additional data A , but only the authentication tag protection of the ciphertext and initialization vector.

Bob can split the received data into parts IV , C and T . This is possible because the lengths of the initialization vector and the authentication tag are parameters of the underlying protocol and must be negotiated beforehand with Alice. Bob can then use the decryption of GCM to validate the authentication tag T . If that fails, the decryption must be aborted altogether because the communication channel must be considered compromised. Else, he can read the additional data A and decrypt the message m .

2.2 Asymmetric encryption

The main problem to solve with symmetric encryption is key distribution. How does Bob get hold of key k in Figure 2.1? Alice and Bob must either meet in person or use a secure communications channel. However, if there already is a secure communication channel, then why not transmit the secret message m through it in the first place? In addition, if n people want to communicate with each other (still, assuming confidential one-to-one communication), then $\frac{n^2-n}{2}$ keys must be exchanged, that is, the number of keys in the system grows quadratically.

In asymmetric encryption (or, equivalently, public-key encryption) the key for encryption is different from the key for decryption. Each user creates a *key pair* of two mathematically linked keys. Data encrypted with one key can only be decrypted using the other key. One key of the pair will be called the *public key* and can be distributed

³Adopted from [3].

freely to anybody who wants to send a message. The other key is the *private key* which must never be disclosed to the public. Figure 2.2 illustrates this process. Alice wants to send an encrypted message to Bob. She asks for Bob's public key E_{Bob} and uses it to encrypt the message. The ciphertext is sent to Bob who can decrypt it using his private key d_{Bob} .

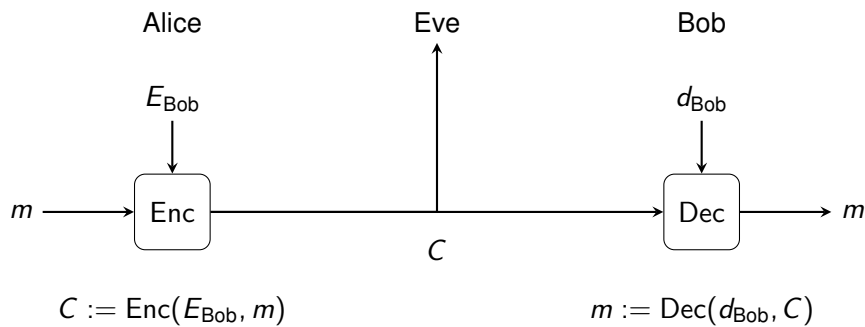


Figure 2.2: Asymmetric encryption and decryption between Alice and Bob.

That is, the following equality holds for all messages m :

$$\text{Dec}(d_{\text{Bob}}, \text{Enc}(E_{\text{Bob}}, m)) = m$$

This approach greatly simplifies the key exchange problem. For a group of n people, only n public keys need to be shared.

So far we only discussed the mathematical link within an asymmetric key pair that guarantees that a ciphertext encrypted with one key can only be decrypted with the other. But we haven't discussed how to construct an asymmetric cipher around that. Without going into detail (which is out of scope for this document), the main idea is to use so-called *trapdoor functions* to construct an asymmetric cipher. The value of a trapdoor function can be easily computed, however, the inverse is hard to impossible to compute without a certain secret piece of information.

The most widely used implementation of such a public-key encryption scheme is RSA [19] which is based on large integer factorization. To derive the private key from a public key, an attacker would have to find the prime factorization of a very large integer N which is the product of two very large prime numbers p and q . There are other asymmetric ciphers which draw their security from other hard mathematical problems like, for example, Elgamal [9] which is based on the discrete logarithm problem. At Data4Life we use the RSA-OAEP encryption scheme [4].

2.3 Hybrid encryption

Symmetric and asymmetric ciphers have certain advantages and disadvantages, some of them being complementary to both [14]. Some properties in which both cipher types differ significantly are data throughput, key sizes and key management effort. Symmetric ciphers are orders of magnitude faster during encryption and

decryption. Also, key sizes are typically smaller. Finally, asymmetric encryption schemes shine when it comes to key management. Because the private key isn't shared, the number of public keys to be shared is much smaller compared to a symmetric key distribution,⁴ and asymmetric key pairs typically have a longer lifetime compared to symmetric keys.

Therefore, if Alice wants to send some large message m to Bob, in practice the hybrid encryption protocol depicted in Figure 2.3 is used: Alice randomly generates a symmetric *session key* k , which she will only use once for the next communication with Bob. Alice encrypts the message m using this session key k , thus producing the ciphertext C_m which can safely be sent to Bob. To allow Bob to decrypt the message, he also needs to securely receive k . This is achieved by asymmetrically encrypting it using Bob's public key, thus producing ciphertext C_k . Bob can decrypt k using his private key and then decrypt the message.

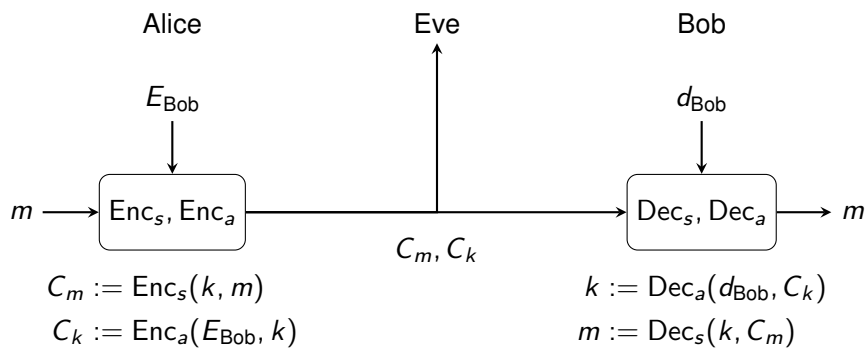


Figure 2.3: Hybrid encryption and decryption between Alice and Bob.

⁴If n persons want to communicate one-to-one, then $\frac{1}{2}(n^2 - n)$ symmetric keys need to be exchanged. Using asymmetric encryption, only n public keys need to be distributed.

3 Secure keys

All healthcare data at Data4Life is end-to-end encrypted. That means, any healthcare data is encrypted on the client before the ciphertexts are sent to the user's data storage. Likewise, when transferring data out of the data storage, ciphertexts are sent to the client where they are decrypted. Furthermore, any communication between client and server is protected via TLS 1.2 [5] currently. TLS 1.3 is planned for the future.

In this chapter we will cover the different cryptographic keys used in our protocols. We use hybrid encryption to encrypt the healthcare data (that is, documents, images and the like) as described in Section 2.3. So far we exclusively considered *healthcare data* or *healthcare documents*. We stick to this simplification for describing the encryption mechanisms we use. In Chapter 4 we will introduce the actual finer-grained data model. However, all concepts discussed below apply there as well.

3.1 Keys overview

When a user first registers her account with Data4Life, she chooses a password w_P which must comply with certain security requirements.¹ Each time at login, after the users was successfully authenticated, the client derives from the password w_P a symmetric key k_P^E . See Section 3.2 below for a precise description of the key derivation process.

Each document is symmetrically encrypted with its own *data key* which is generated by the client when the document is first uploaded. Data keys must be stored alongside their document ciphertext in order to be accessible by multiple clients, for example a browser and a mobile app. Therefore, each data key is symmetrically encrypted using a *common key*. The first such common key² is generated during account creation. A single common key typically protects multiple data keys. Common keys may be shared with a third-party to allow access to specific documents.³ On access

¹It must be at least eight characters total length, contain mixed case letters, contain at least one numeric character (0–9) and at least one special character (such as %, @, or #).

²There might be multiple common keys. This is discussed in Chapter 4.

³Whoever has access to a common key can decrypt all documents whose data keys are encrypted using this common key, provided this party also has access to the document ciphertexts. The server uses business logic to restrict access to only those document ciphertexts that the third party actually was authorized to see.

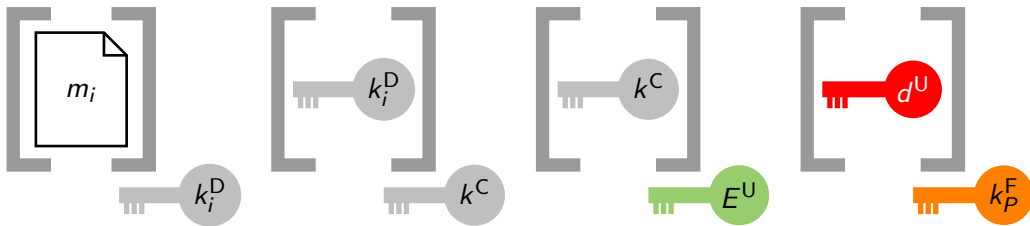


Figure 3.1: Key encryption relationships.

revocation, only the common key needs to be rotated. This entails re-encrypting on the client side all data keys the former common key protects. Common key rotation is a much more efficient process than re-encrypting entire documents with fresh data keys.

Being part of the encryption chain, the common key must also be stored on the server. It is asymmetrically encrypted using the *user public key* which is the public key of a key pair that is also generated during account creation. Lastly, the *user private key* of that key pair is symmetrically encrypted using the password-derived key k_P^F and stored on the server as well.

Let us revisit the key dependencies: Figure 3.1 illustrates the encrypted keys that are stored on the server. We introduce the following short-hand notation to talk more concisely about the various keys in use:⁴ Symmetric keys are denoted by a lowercase k , private keys by a lowercase d and public keys by an uppercase E . The key type is annotated in the superscript: D for data key, C for common key, U for user key, F for a functionally-derived key. In Figure 3.1, gray brackets with an offset key mean that the object in brackets is encrypted using the offset key at the lower right. The document m_i is symmetrically encrypted using its data key k_i^D . The data key is symmetrically encrypted using the common key k^C , which in turn is asymmetrically encrypted with the user public key E^U . The corresponding user private key d^U , finally, is symmetrically encrypted using the password-derived key k_P^F . The key k_P^F is not stored on the server, but derived in the client each time the user logs in with the password w_P .

Let us practice notation and key relationships yet again, but this time from the client's point of view. That is, which steps are to be performed to display a document m_i in the client? Figure 3.2 illustrates this without visual clutter by using the notation established above. We borrow from Figure 3.1 the convention to abbreviate $\text{Enc}(k, m)$ by $[m]_k$. Note how the four ciphertexts in the first row of Figure 3.2 coincide with the ciphertexts depicted in Figure 3.1.

Ultimately, we want to access the plaintext of document m_i . Its ciphertext $[m_i]_{k_i^D}$ is read from the server.⁵ In order to decrypt the document ciphertext, we need

⁴Chapter 7 will introduce the full notation.

⁵We will discuss authorization checks later in Chapter 4. Obviously, the server is not handing out ciphertexts to unauthorized parties.

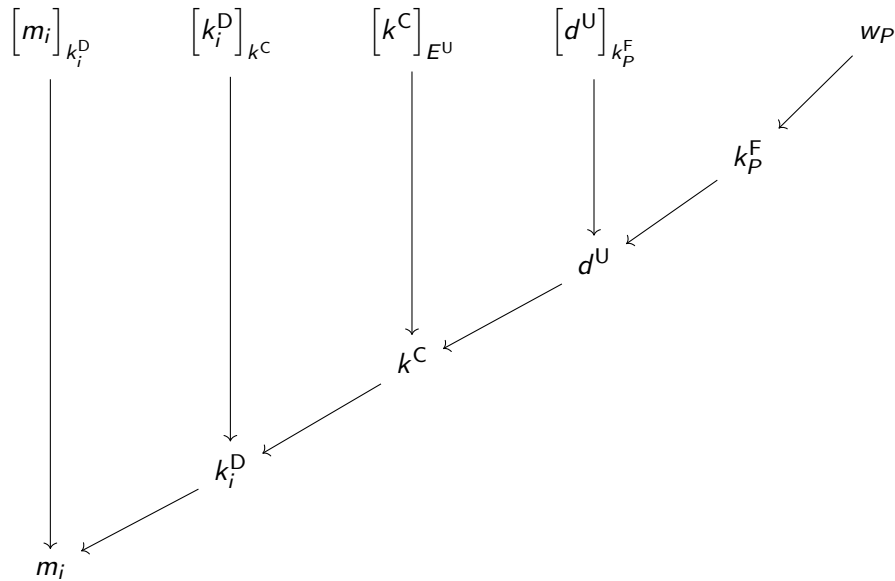


Figure 3.2: Keys and key ciphertexts involved to decrypt document m_i .

Key type	Protected/encrypted by
Data key	Common key
Common key	User public key
User private key	Password-derived key
Password	Remembered by user

Table 3.1: Protection and encryption relationships of the different key types.

the corresponding data key k_i^D whose ciphertext $[k_i^D]_{k^C}$ is read from the server, typically along with $[m_i]_{k_i^D}$. To get access to the data key k_i^D , we need the common key k^C . Its ciphertext $[k^C]_{E^U}$ is also retrieved from the server, and it requires the user private key d^U to be decrypted. The ciphertext $[d^U]_{k_P^F}$ of the user private key is finally decrypted using the password-derived key k_P^F , which exists in the client only ephemerally. The diagonal sequence in Figure 3.2 represents the chain of keys required to decrypt a document. Table 3.1 summarizes the relationships between the key types.

3.2 Account creation

We demonstrated earlier in this chapter that the encryption of a document and its data key requires a couple of other keys to be present in the system, namely the common key, and the public and private user keys.⁶ Those keys are generated during account creation (or, synonymously, account registration).

The *user key pair*, consisting of the user private key and user public key, remains unchanged during a user account's lifetime. All other keys can be changed by the user or are rotated during certain operations. When registering a new account, the user is asked for her e-mail address and a password. A hash of the password is sent to the *haveibeenpwned* API [1] to check whether it might be compromised. If the check passes, the following keys and salts are randomly sampled by the client:

User key pair This is an RSA key pair of length 2048 bits. We denote it by (d^U, E^U) and refer to d^U as the *user private key* and E^U as the *user public key*. This key pair does not change throughout the lifetime of the user account.

Recovery password In case the user forgets her password, it can be reset using this recovery password. It is essentially a second password w_R which is the BIP-39 mnemonic [15] of a random 128 bit number. The user has only one opportunity to download this recovery password, as a PDF file, during the account creation. The respective screen is shown in Figure 3.3. If the user loses this password and forgets her own chosen password, she cannot access the account anymore and all data becomes undecryptable. An example recovery key might look like this:

grunt runway wet horror tent economy
garment photo pause dice achieve soul

Common key The common key consists of 256 bits of randomness that make up an AES key. So far we have denoted it as k^C , but we will add an index like in k_0^C for the remainder of this document, because there may be multiple common keys present in the system simultaneously. The reasons for this are covered in Chapters 4 and 5. However, after account creation, there is initially just a single common key.


Tag encryption key As we will see in Chapter 4, each document can have associated metadata called *tags*. A tag might contain the document type or other, non-healthcare data. We symmetrically encrypt those tags with the *tag encryption key* k^T . This key also remains constant throughout the lifetime of an account.

⁶In reality there is one more key involved: the tag encryption key. It encrypts metadata of the document and we omit it here for brevity. It is explained in Chapter 4.

← Save your recovery key

If you forget your password, the following recovery key lets you reset it.

If you forget your password and lose your recovery key, your account can't be restored, and you can no longer access your data.

 **Your recovery key:**

there poet youth involve month easily print bread spike genre dwarf sail

DOWNLOAD RECOVERY KEY

I understand the importance of my recovery key

I've stored my recovery key in a safe place

Next

Figure 3.3: The user has exactly one opportunity to download the recovery key during account creation.

Data element	Encrypted with
User e-mail address	
User password derivative hash	
Recovery password derivative hash	
User password salt	
Recovery password salt	
User public key	
Encrypted user private key	User-password-derived key
Encrypted user private key	Recovery-password-derived key
Encrypted common key	User public key
Encrypted tag encryption key	Common key

Table 3.2: Data elements sent to the server when creating a new user account.

User password salt and recovery password salt Two salt values with 16 bits of randomness each are generated for the user and the recovery password. The user salt is used to derive the symmetric key k_P^F . The recovery salt is used for a similar reason, in case the user password is to be reset (see Section 3.4 below).

User password derivative and recovery password derivative The client does never send raw passwords to the server during account creation or login, but rather a key derivative of it. We use the output of the PBKDF2 [13] key derivation function (10,000 iterations with zero salt and SHA-256) as a derivative value.

Clicking “Next” on the screen of Figure 3.3 will send the data given in Table 3.2 to the server where it makes up the user’s account. Before password derivatives are sent to the server for storage, they are hashed using *bcrypt* [18] with a cost parameter of 10. The resulting hashes are encrypted at rest with a secret server key using AES in GCM.

3.3 Login

Apart from authentication and authorization, the login procedure serves another important purpose called *client approval*. This approval establishes fresh keys for encrypting and decrypting documents. Any data-related server API endpoint requires a valid JSON Web Token (JWT [12]) to be contained in client requests. Login and client approval, if successful, will result in such a JWT.

3.3.1 Authentication and authorization

To log in, the user is asked for her e-mail address and password. The e-mail address and the password derivative (see Table 3.2) are sent to the server. If the e-mail is known, the hash of the password derivative matches the one stored on the server, and a second factor (in our case a PIN received via text message on the phone) was

provided, the user is considered authenticated. The login request will also contain a list of so-called *scopes*, which are strings designating specific access rights.⁷ If the scopes are allowed for the user and the requesting application (that is, the application that triggered the login), authorization succeeds, and client approval is commenced.

3.3.2 Client approval

As described above, the user key pair is a long-lived key that essentially protects all other keys in the system (see Figure 3.1, where it is depicted at the end of the “key chain”).⁸ Technically, the client could now exercise the flow depicted in Figure 3.2 to decrypt documents of the user, and also to encrypt new documents when proceeding in the opposite direction. However, this requires the user private key to be present in the client for as long as the login session lasts. In order to reduce attack surface, we decided to require every login session to provide its own ephemeral so-called *application key pair*.

Prior to commencing the login process, the client generates a fresh asymmetric application key pair (d_s^A, E_s^A) . (Let the index s denote the word *session*.) If authentication and authorization succeed, the server sends the ciphertexts of the user private key and of all common keys to the client:

$$\left\{ \left[d^U \right]_{k_P^F}, \left[k_0^C \right]_{E^U}, \left[k_1^C \right]_{E^U}, \dots \right\}$$

These ciphertexts are decrypted, re-encrypted using the application public key, and sent back to the server which stores them:

$$\left\{ E_s^A, \left[k_0^C \right]_{E_s^A}, \left[k_1^C \right]_{E_s^A}, \dots \right\} \quad (3.1)$$

The client can decrypt any documents that are currently stored in the user account, regardless which client uploaded them (and, hence, which common keys were used to protect the data keys). This is because it has access to all common keys that are in the user account at the moment. The client approval is now complete. The server will return a JWT that the client must include with all subsequent server calls to prove authenticity and authorization.

3.4 Account recovery

In Section 3.2 we briefly mentioned that the recovery password, which is created and downloadable during account creation, can be used to get access to the account in case the user has forgotten her password w_P . Without password w_P , no key k_P^F can

⁷It would not be wrong to think of such a scope string as a user role. The notion of a scope is adopted from the OAuth2 authentication protocol [11].

⁸You might argue that the user-password-derived key is the most crucial one—and that is correct—however, the password (hopefully) persists exclusively in the user’s memory, and the derived key only ephemerally in the client.

be derived, which is required to start the “decryption chain” depicted in Figure 3.2. However, along with $[d^U]_{k_P^F}$ we also generated $[d^U]_{k_R^F}$ and stored it on the server during account creation. By choosing a dedicated recovery flow during login and by entering the recovery password w_R , the client can decrypt the user private key d^U , ask the user for a new password $w_{P'}$ and send the new ciphertext $[d^U]_{k_{P'}^F}$ back to the server.

3.5 Key maps

So far we have discussed a number of key topics. Before we continue, let us establish a common understanding of how the server stores key ciphertexts for each user. Figure 3.4 shows an example. At any given time, a user’s data store contains the following individual key ciphertexts:

- User private key encrypted with the user password-derived key: $[d^U]_{k_P^F}$
- User private key encrypted with the recovery password-derived key: $[d^U]_{k_R^F}$
- Tag encryption key encrypted with the current common key k_x^C : $[k^T]_{k_x^C}$

Figure 3.4 depicts these three ciphertexts in the first row (the current common key is k_1^C). Further, each user account contains one or more key maps (or, more precisely, common key maps). A key map contains for each asymmetric key pair (user key pair or application key pairs) the corresponding public key and the ciphertexts of all relevant common keys encrypted using the respective public key.⁹ One common key ciphertext in each key map is marked as the key map’s *current common key*. We denote this by a little dot above the corresponding ciphertext in Figure 3.4. Finally, a key map may also contain the ciphertext of the tag encryption key.

There is at least one key map π_\emptyset present in a user account. It holds the user public key and ciphertexts of all the common keys stored in the system for the user. When a user logs in, the login session that is established, results in a new key map, like π_s in Figure 3.4. The contents of that key map are the result of the client approval (see Equation 3.1 above). For example, π_s contains the session public key E_s^A and the ciphertexts of the three common keys k_0^C , k_1^C and k_h^C that were present at the time of client approval. Key map π_h was not created during a login client approval, but during the so-called onboarding of a third party h (see Chapter 5 for details). For now it is sufficient to accept that the key maps may contain ciphertexts of several common keys. There is exactly one dot above one common key in each map, denoting the current common key for the respective session or third party.

⁹Relevant means that not every key map contains the ciphertexts of all common keys.

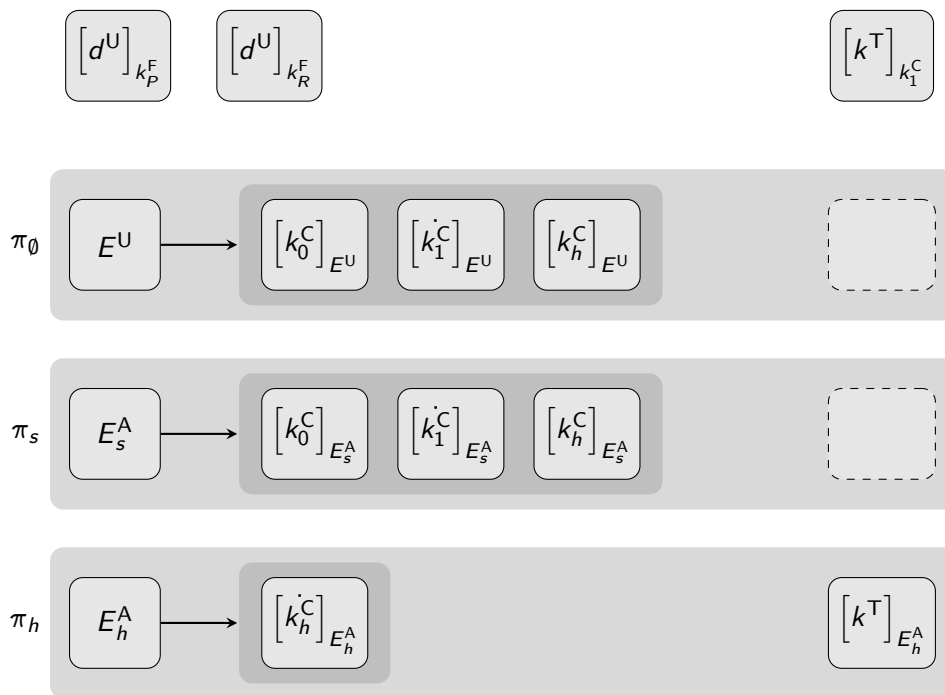


Figure 3.4: Example of a user account with three key maps containing ciphertexts of different common keys. Key map π_0 is present in every account and maps the user public key to the ciphertexts of all common keys known to the account. Any other key map may contain a subset of them. Key map π_s contains the ciphertexts of the same three common keys, while key map π_h contains only the ciphertext of common key k_h^C . π_h also contains the ciphertext of the tag encryption key.

4 Data model and data access

This chapter describes the full data model of the data storage together with the cryptographic protocols used to write and read data by a user. We discussed the main ideas in Chapter 3, but we left out some details covered here.

So far we talked about documents that users can maintain in their data storage, and that is certainly a metaphor worth adhering to. The actual data model is more granular, though. There is the notion of a *record*, which can be considered a folder that contains so-called *attachments*. When the user selects one or more files from disk to be uploaded, a new record is created and the files are stored as associated attachments. It is not a mistake to consider the unit of a record and its attachments as a document in the sense we used it so far. Each record can have zero or more associated *tags*. Tags are strings (typically key-value pairs) which can contain operational data such as a data type. The record itself has a body which is a FHIR resource [6] and as such essentially just a JSON string.

Figure 4.1 illustrates the relationships. The record r_i be a JSON FHIR resource of type DocumentReference. The three attachments $a_{i,1}$, $a_{i,2}$ and $a_{i,3}$ could be PDF files, for example. The two tags $t_{i,1}$ and $t_{i,2}$ be some metadata (for example, tag $t_{i,1}$ could be `uploadedVia=mobile`). The attachment IDs, here, exemplary, 1 and 2, but in reality these would be random UUIDs, are stored inside the record's FHIR body.

The record body is symmetrically encrypted using its corresponding data key k_i^D , which is used for this record only and nowhere else. In case a record is updated, a new data key is generated and used.¹ Each record is associated with another symmetric key, the *attachment key* k_i^N , which is used to encrypt all its attachments. Each tag is symmetrically encrypted using the account-wide *tag encryption key* k^T . As shown in Table 3.1 in Chapter 3 (page 16), the data key (and also the attachment key) are encrypted using the common key, or, more precisely, the *current* common key k_x^C . For technical reasons we tolerate multiple common keys to be in the system. However, at any given time there is only one marked as current key and used for encrypting new data uploaded by the user. Table 4.1 summarizes the data that is stored on the server for our example record r_i .

Record and attachments are encrypted using AES-256 in GCM.² The initialization vector consists of 12 bytes of randomness and the authentication tag is 16 bytes long.

¹Essentially, a record update is deletion followed by creation of a new record.

²See Section 2.1.6 for details.

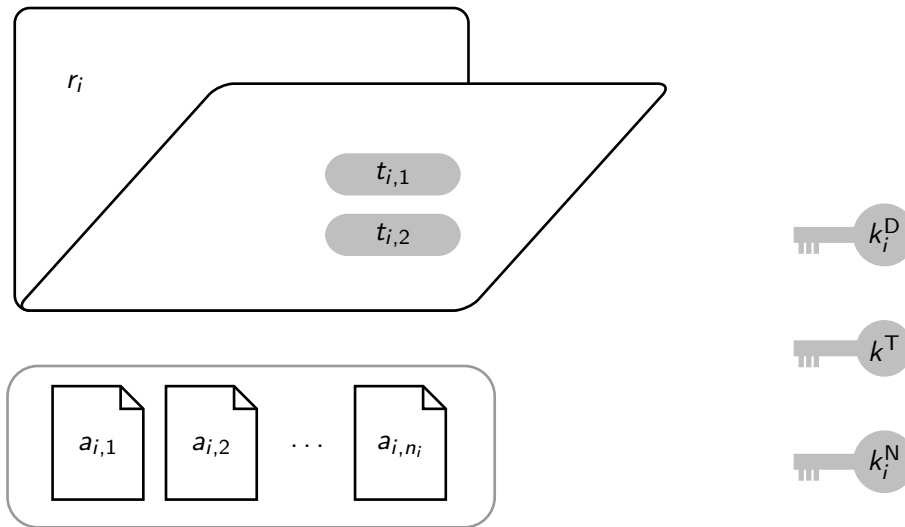


Figure 4.1: Data model example.

Data element	Formal notation
Encrypted record body	$[r_i]_{k_i^D}$
Encrypted data key	$[k_i^D]_{k_x^C}$
Encrypted attachment key	$[k_i^N]_{k_x^C}$
Encrypted tags	$[t_{i,1}]_{k_i^T}, [t_{i,2}]_{k_i^T}$
Encrypted attachments	$[a_{i,1}]_{k_i^N}, [a_{i,2}]_{k_i^N}, [a_{i,3}]_{k_i^N}$
Common key ID	x

Table 4.1: Data elements stored for a record r_i with two tags and three attachments.

The binary payload that gets stored for an encrypted record or encrypted attachment consists of the concatenation of the initialization vector, the actual ciphertext, and the authentication tag as depicted in Figure 4.2. Tags are symmetrically encrypted using AES-256 in CBC mode with the initialization vector consisting of 16 zero bytes.

IV (12 byte)	Ciphertext (n byte)	Auth'n Tag (16 byte)
-----------------	---------------------------	-------------------------

Figure 4.2: Data layout for symmetrically encrypted records and attachments.

Let us now revisit which steps are taken when a user wants to display the attachment $a_{i,1}$ of record r_i . Since the attachment ID is stored in the record body, it must be decrypted first. This is shown in Figure 4.3 (this chart is essentially identical to Figure 3.2 up to the namings of some artifacts, like message m_i versus record r_i and the like). After decrypting the record r_i , the corresponding attachment IDs are known. This allows to retrieve the ciphertext $[a_{i,1}]_{k_i^N}$ and proceed analogously for decrypting it as depicted in Figure 4.4.

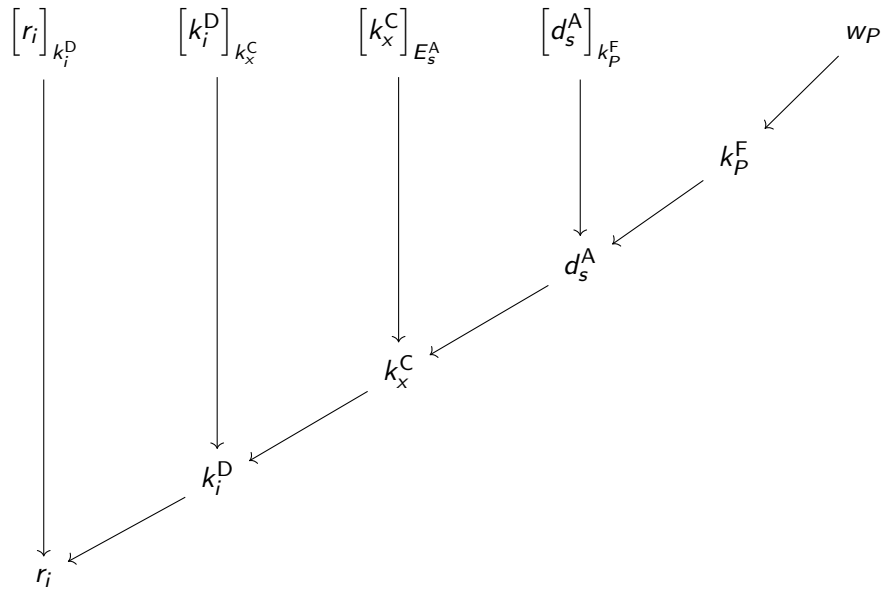


Figure 4.3: Data decryption flow for reading record r_i .

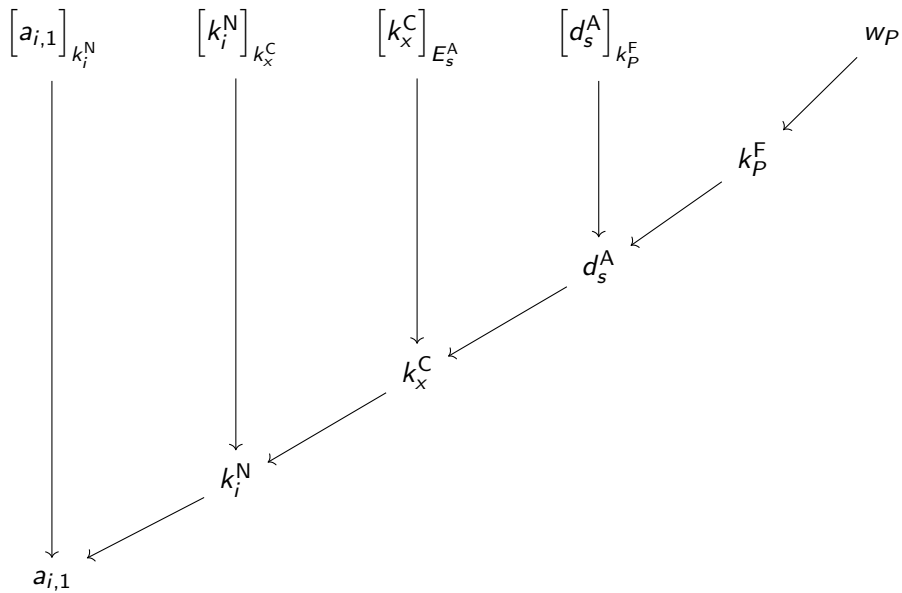


Figure 4.4: Data decryption flow for reading attachment $a_{i,1}$ after having read the record r_i .

5 Data ingestion

One way of adding documents to a user's data storage is to use the web application or mobile application and upload data manually. This was covered in Chapter 4. We provide another avenue by which healthcare documents generated by hospitals (such as doctor letters or discharge letters) can be transferred securely directly into the user's data storage.

The main challenge is to enable a third party, for example a hospital, to securely add documents to a user's data storage while not allowing them any read access. As we saw in Chapter 3, a common key protects data keys, that is, data keys are encrypted by a common key. Any party who holds a common key can decrypt all documents¹ whose data keys are encrypted with it. Note, that this argument implies that a third party would have access to all ciphertexts of those documents, which, in practice, will not be the case, of course.

The main idea for enabling a third party to add documents is to generate an additional common key which is securely transferred to the third party for encrypting the data keys of the new documents. The upload of encrypted documents by a third party is cryptographically identical to the upload of a document by the user, except that the common key used for data encryption keys is different. Table 5.1 illustrates this subtle but important difference. When a user uploads a document, its data key is encrypted using the current common key (here assumed to be still the initial one k_0^C). When a hospital h uploads a document for the user, the data key is encrypted using the hospital's dedicated common key k_h^C . Note that the hospital would hold such a dedicated common key for each user who granted them write access. That is, we should more precisely denote the hospital's common key of user u as $^{(u)}k_h^C$ to stress that user binding. In the remainder, however, we assume the user to be implicitly fixed and omit the (u) decoration.

5.1 Onboarding

Before a third party can add documents to a user account, it has to undergo onboarding in a special type of login process. The third party h generates an application key

¹For the sake of simplicity, we again talk about documents, but do, of course, recall, that they consist of records and attachments (see Chapter 4). Nevertheless, all arguments and ideas put forth w. r. t. documents in this chapter carry over to records and attachments.

Party	Common key	Ciphertexts sent for document m_j
User	k_0^C	$[m_j]_{k_j^D}$, $[k_j^D]_{k_0^C}$, $[t_{j,*}]_{k^T}$
Hospital h	k_h^C	$[m_j]_{k_j^D}$, $[k_j^D]_{k_h^C}$, $[t_{j,*}]_{k^T}$

Table 5.1: Ciphertexts sent to the server by user vs. a third party such as a hospital.

pair (d_h^A, E_h^A) and sends the public key E_h^A to the authentication and login application at Data4Life.² The application detects that we are in an onboarding flow instead of a regular login flow, and performs the following steps:

1. Generate a new common key k_h^C for third party h .
2. Encrypt it using the provided public key: $[k_h^C]_{E_h^A}$
3. Execute a modified client approval:
 - Add to each non-foreign key map π_x in the user account the ciphertext of the new common key: $[k_h^C]_{E_x^A}$. *Non-foreign* refers to key maps which do not hold keys for another third party. In other words, do not add the ciphertext to any key map which was established using this very protocol.
 - Add the new key map $\pi_h = \left(E_h^A, \{ [k_h^C]_{E_h^A} \} \right)$ to the user account.
 - Add to the key map π_h the tag encryption key ciphertext $[k^T]_{E_h^A}$.
4. Allow third party h to read the contents of π_h , that is, the ciphertexts $[k_h^C]_{E_h^A}$ and $[k^T]_{E_h^A}$, at any later point in time by issuing a JWT that allows requests against the server. See Section 3.3 for details.³

This concludes the onboarding process of a third party.

5.2 Document upload

Adding a document by the third party works as follows:

²The very same web application that implements the account creation, login, client approval and account recovery flows.

³In fact, we are issuing two tokens: a short-lived JWT for immediate use and a long-lived refresh token, which can be used at a later point in time to request a new JWT (this is identical to the principles of OAuth2 [11]). Since the third party will rarely add documents immediately after onboarding, the JWT will likely expire. It is the refresh token that is kept by the third party and exchanged for a fresh JWT prior to the actual document upload.

1. Exchange refresh token for a JWT to authenticate subsequent server requests.
2. Request the ciphertexts of π_h from the server: $[k_h^C]_{E_h^A}$ and $[k^T]_{E_h^A}$
3. Decrypt k_h^C and k^T using the private key d_h^A
4. Encrypt a document m_j as usual, that is, produce: $[m_j]_{k_j^D}$ and $[k_j^D]_{k_h^C}$
5. Encrypt the necessary tags, say, $t_{j,1}$ and $t_{j,2}$ with the tag encryption key k^T .
6. Upload ciphertexts $[m_j]_{k_j^D}$, $[k_j^D]_{k_h^C}$, $[t_{j,1}]_{k^T}$ and $[t_{j,2}]_{k^T}$ to the server.

Even though the third party did add a document, there is no way for them to decrypt any document ciphertext other than the ones they uploaded. Assume the ciphertext of some document m_x leaks. The corresponding data key k_x^D is encrypted with a common key different from k_h^C and thus the third party cannot decrypt it. There is no way for party h to decrypt any other common key. Its key map does not include them and other key maps, should the party get hold of them via an attack, are encrypted with public keys other than E_h^A .

5.3 Document access by client

Let us verify that the user can read the uploaded documents using any approved client (that is, any application for which there exists a key map in the account). We can, in fact, reuse Figure 3.4 on page 22 to illustrate this. The key map π_h was created during the onboarding of third party h . Let the key map π_s belong to a session that the user just logged in to. Assume further, that party h has uploaded a doctor's letter m_j . In order to decrypt it in the client (session s), the following steps are performed:

1. Download ciphertexts $[m_j]_{k_j^D}$ and $[k_j^D]_{k_h^C}$.
2. Key map π_s contains $[k_h^C]_{E_s^A}$. Decrypt it using d_s^A to get k_h^C .
3. Use k_h^C to get k_j^D , which then is used to get m_j .
4. Tags $t_{j,1}$ and $t_{j,2}$ can be decrypted, because k^T is decryptable via k_1^C which is in key map π_s .

The client, which can only access its session key map π_s , can decrypt the data key because the required common key k_h^C was inserted into π_s during approval of the third party onboarding.

6 Data donation

Users of the PHDP can choose to donate (parts of) their data for research purposes. The research tasks are carried out on the Analytics Platform (ALP). The data is decrypted on the client, anonymized, re-encrypted, and then transferred into a special storage where it awaits access by ALP. In this chapter we introduce the requirements for data donation and discuss the final design.

6.1 Requirements

Let us assume that a user with ID u wants to donate the document m_i . Even though the ID u will be a random UUID, let us assume for the sake of argument, that it is a more recognizable value like her e-mail address. Similar to other chapters, we talk about documents for simplicity, but do recall that they consist of records and attachments (see Chapter 4). All arguments and ideas in this chapter, of course, carry over to records and attachments. Last, we assume that the document m_i was anonymized by some means.

De-identification One requirement is that donated data must be completely detached from the users' identities. That is, it must not be possible to correlate the donated document m_i with the user ID u .

Data quality Users may regularly donate documents (that is, not in one batch). In order to be of value, on ALP side it must be known which documents came from the same donors.¹ Effectively, each user is assigned another ID u^{ALP} (referred to as ALP ID), under which all donated documents of that user are collected by ALP. It must not be possible to infer u from u^{ALP} .

Consent We need to make sure that at the time of donation users had given consent to donate their data. In case of a user claiming that we store data without her consent, we must be able to prove her consent. The consent is bound to the user ID u , because the donated data is not traceable to the donor anymore.² The donation workflows must check for a valid consent before accepting data donations.

¹Otherwise, longitudinal analyses would not be possible.

²Also, it is anonymized, which means it is no personal data anymore.

No impersonation Users must not be able to donate data such that it gets associated under another user's ALP ID.

6.2 Design Decisions

In order to understand the rationale behind the sequence diagrams in Figures 6.1 and 6.2, we now briefly cover the steps in reverse order, before discussing them in full detail in Sections 6.3 and 6.4.

After the document m_i was donated by user u , the ALP must be able to collect the ciphertext of m_i from a special storage where the donation service deposited it. The document is encrypted with a special data key k_i^{Don} , which itself gets encrypted with the public service key of the ALP $E_{\text{ALP}}^{\text{S}}$. Finally, this data must be associated with an ALP ID u^{ALP} which is unique and different for each user u . That is, the following structure must arrive at the donation service to conclude a successful donation of document m_i by user u :³

$$\left\langle u^{\text{ALP}}, \left[m_i \right]_{k_i^{\text{Don}}}, \left[k_i^{\text{Don}} \right]_{E_{\text{ALP}}^{\text{S}}} \right\rangle$$

Even though the document m_i is anonymized (and thus does not anymore constitute personal information), we encrypt it for reasons that become clear below. One challenge is to maintain for each registered donor u an ALP ID without being able to invert the mapping and resolve the user ID u given an ALP ID u^{ALP} . We solve this problem by introducing for each donor a donation key pair $\left({}^{(u)}d^{\text{Don}}, {}^{(u)}E^{\text{Don}} \right)$ which itself is securely stored in the users PHDP. The public key⁴ E^{Don} serves two purposes. First, the donation service maintains a mapping $E^{\text{Don}} \mapsto u^{\text{ALP}}$, which treats E^{Don} as a proxy for the user ID u , without, however, being able to reverse-map it to u . Second, we are using d^{Don} to sign the ciphertext of the donated document, which allows the donation service to verify the data provenance (again, without being able to identify the user u). We will now revisit the registration and donation processes in full detail.

6.3 Register as a donor

Users must register as donors before they can donate documents. During the registration a user must give consent to the data donation terms and conditions. Further, a donation key pair is created for the user, the public key of which will be assigned to a newly created ALP ID. The sequence diagram in Figure 6.1 illustrates all steps of a successful registration. Error handling is omitted for brevity, but it is discussed in the narrative below. In previous chapters we referred to the PHDP and

³We use the angle bracket notation $\langle \cdot \rangle$ to denote a structure of data elements. Think of it as a tuple or some JSON object.

⁴We omit the superscript (u) in the remainder, because it is clear that we are discussing everything in the context of user u .

its services simply as *the server*. The donation solution introduces two new services, consent service and donation service, which would naturally be covered under the umbrella of PHDP, and thus, *the server*. Because they are crucial to the donation solution, they are depicted alongside the PHDP in the sequence diagrams. The two aforementioned services and the ALP maintain service key pairs that are used for either signing data or encrypting requests. The public keys of the key pairs are implicitly assumed to be known to the client (that is, we do not explicitly draw the key retrieval requests in the sequence diagrams to save space) and all services. The key pairs are denoted as follows:

$$\begin{aligned} \text{Consent service (CS):} & \quad \left(d_{\text{ConSig}}^S, E_{\text{ConSig}}^S \right) \\ \text{Donation service (DS):} & \quad \left(d_{\text{Don}}^S, E_{\text{Don}}^S \right) \text{ and } \left(d_{\text{DonSig}}^S, E_{\text{DonSig}}^S \right) \\ \text{Analytics Platform (ALP):} & \quad \left(d_{\text{ALP}}^S, E_{\text{ALP}}^S \right) \end{aligned}$$

6.3.1 Steps 1 – 7: Request consent

If a user want to become a data donor **1**, she must give her consent to the current data donation terms and conditions. The client requests the latest version of the consent narrative ξ in step **2**. The consent service caters for different types of user consents, not only for data donation. The value ξ is a placeholder for a specific type of consent narrative which exists in different versions, the latest being version V . Hence, the consent service returns the latest narrative of the donation consent T_V^ξ in step **3**. The user approves the consent (**4**) and this fact together with the current timestamp τ is sent to the consent service (**5**) which stores it in a consent database (**6**, **7**). If the user rejects the consent, the donor registration process is aborted.

6.3.2 Steps 8 – 10: Create donor key pair

The client now generates a new asymmetric key pair $(d^{\text{Don}}, E^{\text{Don}})$ in step **8** and stores it in the user's PHDP records (**9**, **10**). This key pair is used for each subsequent document donation.

6.3.3 Steps 11 – 18: Sign registration request

As mentioned in Section 6.2 above, the goal is to let the donation service generate a new ALP ID and associate it with the key E^{Don} (or detect that the user is already registered as a donor and abort the registration). Also, the donation service must make sure that the key E^{Don} belongs to an existing user (and not some attacker trying to spam the database) who has given the correct type of consent.

The donation service API cannot be authenticated using the PHDP access token (JWT) that was discussed in Section 3.3.2, because it contains—among other things—the user ID u , which must never reach the donation service in any way. The donation

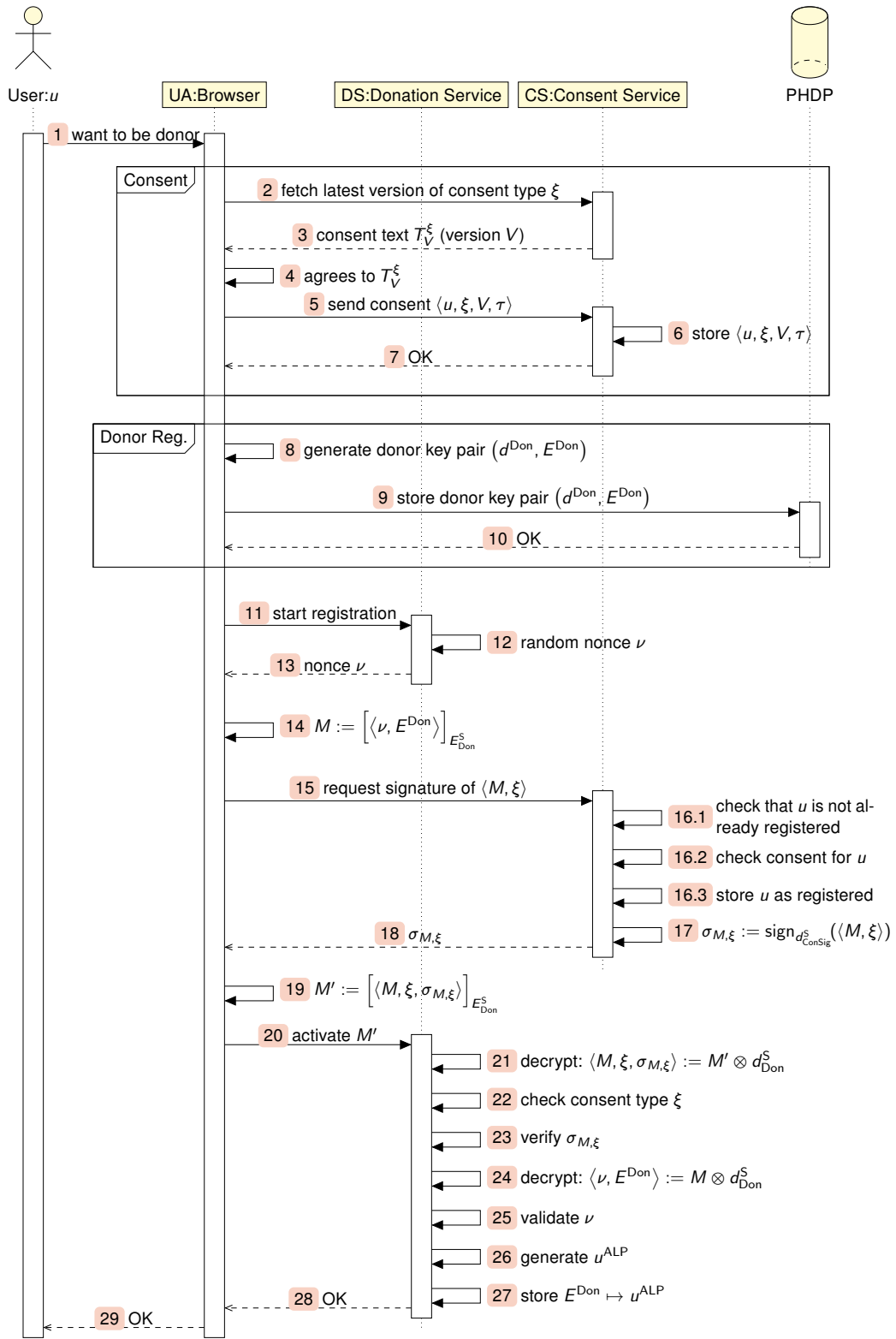


Figure 6.1: User registers as donor.

service issues a random nonce ν instead (11 – 13), which acts as a short-lived session ID for the remainder of the registration process.

The client now creates the message M by encrypting both the users donation public key and the nonce with the donation service public key (14):

$$M := \left[\left\langle \nu, E^{\text{Don}} \right\rangle \right]_{E_{\text{Don}}^{\text{S}}}$$

If we sent this ciphertext to the donation service already now, it could decrypt it using $d_{\text{Don}}^{\text{S}}$ and proceed checking for an existing ALP ID for E^{Don} (in which case we would abort the registration process, otherwise proceed and associate a new ALP ID with E^{Don}). However, there would be no way to verify whether the user gave her consent. Remember, the key E^{Don} cannot (and must not) be back-traceable to the user ID u . However, consent is given and recorded only under such user IDs (and never under a donation public key).

This is why the ciphertext M is sent to the consent service (15) which is requested to provide evidence (in form of a signature) that user u agreed to the consent narrative of type ξ ⁵. Step 16.1 checks whether the user u is not already registered as a donor. If the lookup in the consent database for user u and consent type ξ succeeds (16.2), user u is remembered as registered⁶ (16.3), and the pair $\langle M, \xi \rangle$ is signed using the consent service private key $d_{\text{ConSig}}^{\text{S}}$ (17). The resulting signature $\sigma_{M,\xi}$ is returned to the client (18).

6.3.4 Steps 19 – 29 : Activate registration

Now that we have a signed version of the donation public key E^{Don} , we could send it to the donation service, which in turn could verify that the user (whose real identity u the service will never know) did, in fact, give the right consent. This would, however, entail sending M and ξ over the network for a second time (it did already leave the client in step 15). An attacker knowing the registration protocol (which is not secret) could detect this and learn important information from it. This is why step 19 encrypts the data (now also including the signature $\sigma_{M,\xi}$ a second time with the donation service key $E_{\text{Don}}^{\text{S}}$, producing ciphertext M' :

$$M' := \left[\left\langle M, \xi, \sigma_{M,\xi} \right\rangle \right]_{E_{\text{Don}}^{\text{S}}} = \left[\left\langle \left[\left\langle \nu, E^{\text{Don}} \right\rangle \right]_{E_{\text{Don}}^{\text{S}}}, \xi, \sigma_{M,\xi} \right\rangle \right]_{E_{\text{Don}}^{\text{S}}}$$

This ciphertext is sent to the donation service (20) which can decrypt the following data elements (steps 21 and 24):⁷

$$\nu, E^{\text{Don}}, \xi, \sigma_{M,\xi}$$

The donation service carries out the following verifications (the order does not matter, all checks must pass):

⁵As mentioned earlier, think of ξ as the narrative type “donation consent”.

⁶A new registration is possible only after waiting period has passed after a deletion proof was received. See Section 6.5 below.

⁷We denote decryption of ciphertext M using key k by $M \otimes k$.

- 22 Check that the consent is of correct type ξ .
- 23 Verify signature $\sigma_{M,\xi}$ using E_{ConSig}^S against $\langle M, \xi \rangle$. That is, verify that consent was given.
- 25 Check whether the nonce ν is, in fact, one that was issued by the donation service recently.

If all checks pass, the donation service can generate a new ALP ID u^{ALP} (26) and store it for key E^{Don} (27). This concludes the donor registration process (steps 28 and 29 being just formal returns).

6.4 Donate a document

After users have given consent and are registered as a donor, they can donate documents. We illustrate the donation of a document m_i by user u in Figure 6.2. As mentioned above, we assume the document m_i was anonymized in the client.

6.4.1 Steps 1 – 12 : Prepare donation

The user selects document m_i for donation (1), after which the client retrieves the donation key pair $(d^{\text{Don}}, E^{\text{Don}})$ from PHDP (steps 2 and 3). Note, that this is simplified, because data access follows the protocol described in Chapter 4, that is, we would read the ciphertexts of records, attachments and data keys, and decrypt them using the common key.

Similar to steps 11 – 13 in Figure 6.1, we request a random nonce ν (steps 4 – 6) that acts as a short-lived session ID for the subsequent communication with the donation service.

Steps 7 – 12 are identical to steps 14 – 19 in Figure 6.1.

6.4.2 Steps 13 – 15 : Encrypt and sign document

The client generates a symmetric key k_i^{Don} in step 13 which is used to encrypt the document m_i in step 14, producing the ciphertext C_i . This ciphertext is signed using the user's private donation key d^{Don} (15) which results in the signature σ_{C_i} . This signature is later verified by the donation service.

6.4.3 Steps 16 – 26 : Verify and store donated document

The following data elements are sent to the donation service in step 16:

$$\left[\left\langle \left[\left\langle \nu, E^{\text{Don}} \right\rangle \right]_{E_{\text{Don}}^S}, \xi, \sigma_{M,\xi} \right\rangle \right]_{E_{\text{Don}}^S}, [m_i]_{k_i^{\text{Don}}}, [k_i^{\text{Don}}]_{E_{\text{ALP}}^S}, \sigma_{C_i}$$

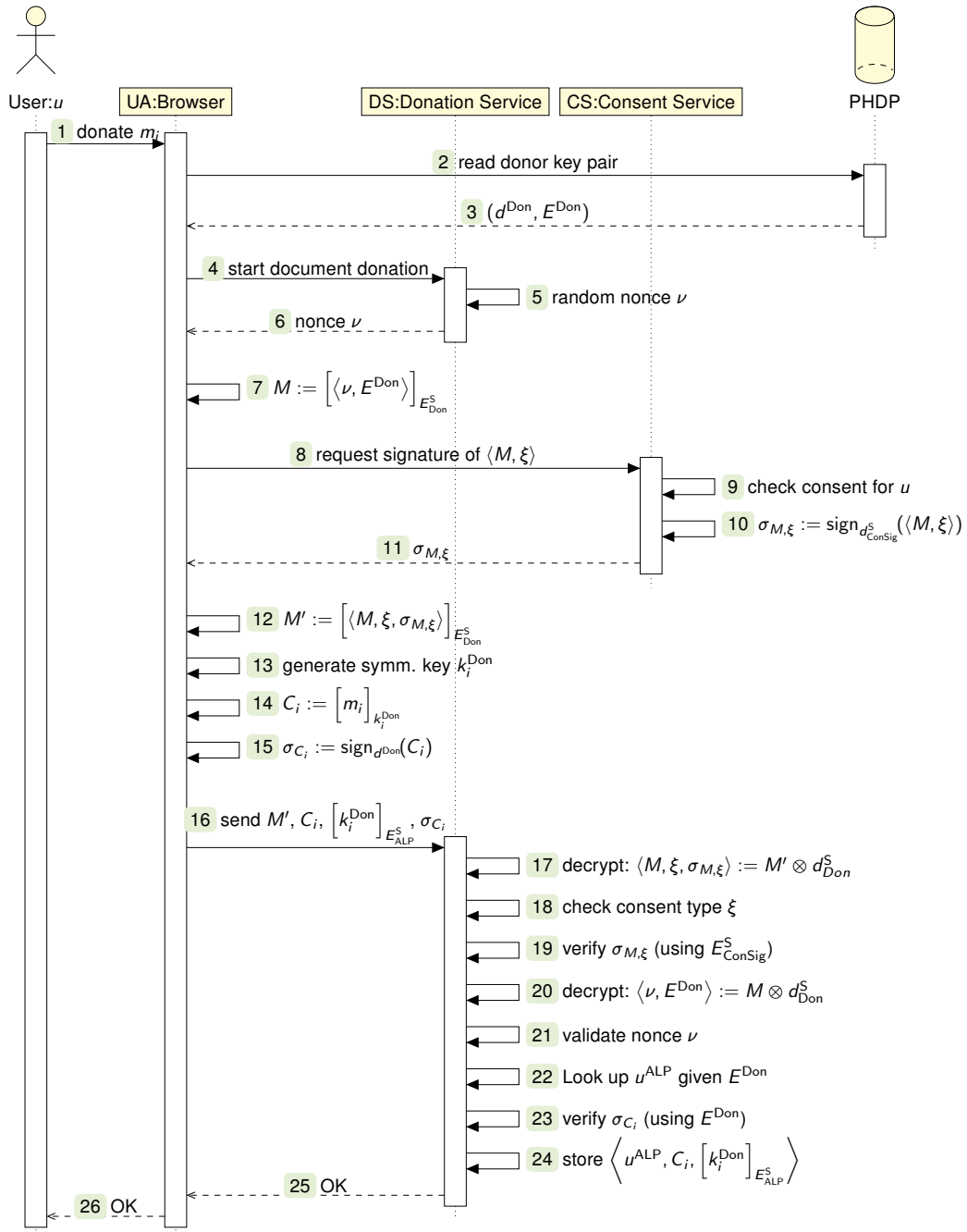


Figure 6.2: User donates document m_i .

This enables the donation service to decrypt and assemble the following data elements (steps 17 and 20):⁸

$$\nu, E^{\text{Don}}, \xi, \sigma_{M,\xi}, \sigma_{C_i}, C_i, [k_i^{\text{Don}}]_{E_{\text{ALP}}^{\text{S}}}$$

The donation service carries out the following verifications (the order does not matter, all checks must pass):

- 18 Check that the consent is of correct type ξ .
- 19 Verify signature $\sigma_{M,\xi}$ using $E_{\text{ConSig}}^{\text{S}}$ against $\langle M, \xi \rangle$. That is, verify that consent was given.
- 21 Check whether the nonce ν is, in fact, one that was issued by the donation service recently.
- 22 Check that there is an ALP ID for key E^{Don} and retrieve it as u^{ALP} .
- 23 Verify signature σ_{C_i} using E^{Don} against C_i .

If all checks pass, the ciphertexts of the donated document can be stored under the ALP ID:

$$\left\langle u^{\text{ALP}}, [m_i]_{k_i^{\text{Don}}}, [k_i^{\text{Don}}]_{E_{\text{ALP}}^{\text{S}}} \right\rangle$$

This corresponds to the objective we discussed in Section 6.2. This concludes the document donation process (steps 25 and 26 being just formal returns). Note that the donated documents are now stored at the donation service. In order to analyze them in the ALP they must be transferred to it. This transfer does not happen immediately. Otherwise, one could correlate traffic between the donation service, ALP and the other services to possibly link donation requests. Each donated document is assigned a retention date before which it will not be transferred to the ALP. The retention date is chosen such that documents arriving at different dates will get assigned identical retention dates. This leads to batches of documents that can be transferred to the ALP. As a further security means the set of ALP IDs of the documents awaiting transfer must exceed a minimum size before they are transferred.

6.5 Revocation

Users can at any time revoke their consent to be a data donor. Revocation comes in two forms: temporary and permanent. Temporary revocation suspends the ability of the user to donate data, but it keeps the donation key pair. That is, when the user decides to be a donor again, the data can be associated with the same ALP ID it used to. Permanent revocation deletes the user's donor key pair. That is, when the user wants do become a donor again, she would need to register anew as if she were never a donor, thus creating a new donor key pair. Newly donated data ends up at the ALP under a new ALP ID and is seen as if coming from a different person. We illustrate the revocation steps in Figure 6.3.

⁸As before, we denote decryption of ciphertext M using key k by $M \otimes k$.

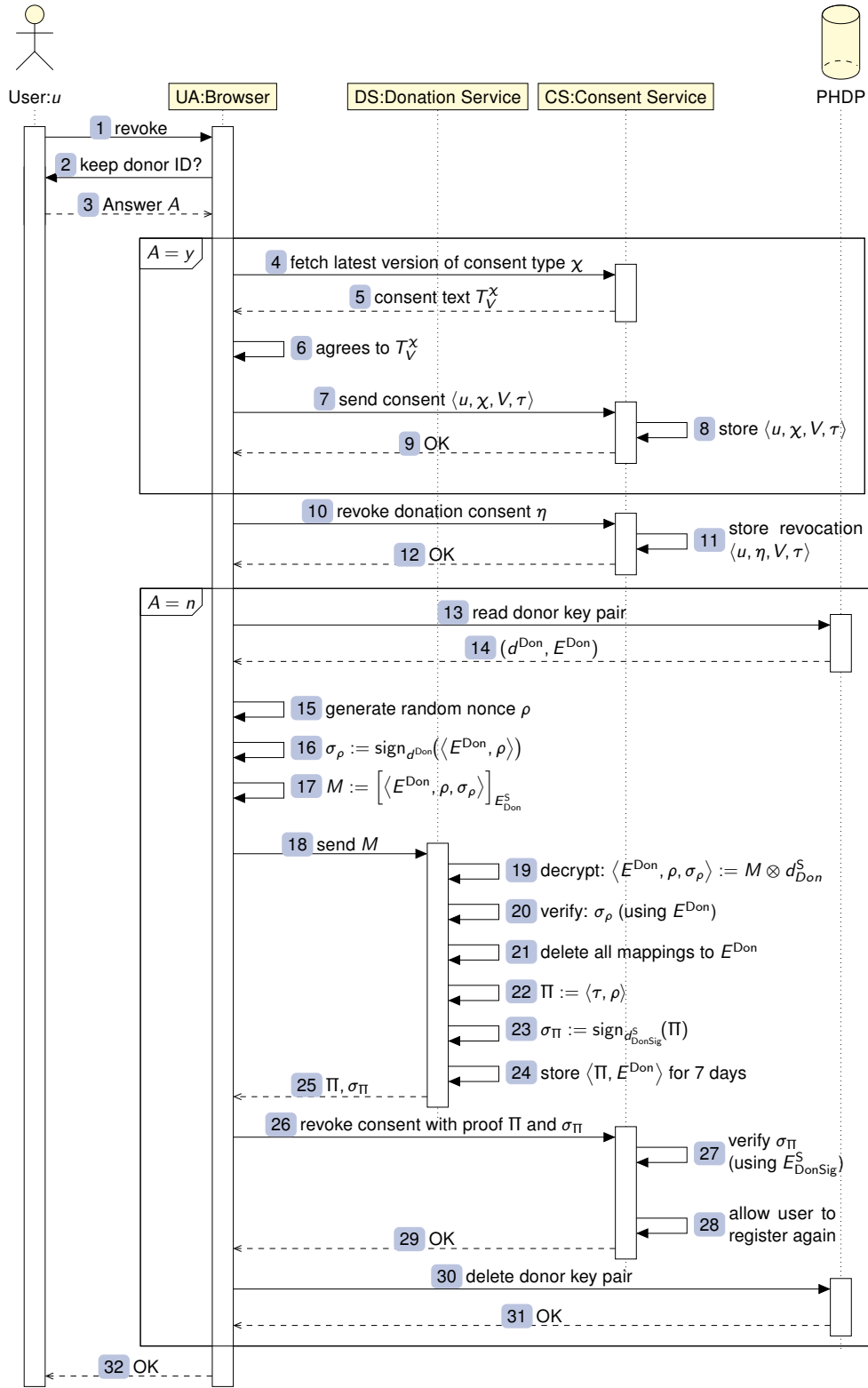


Figure 6.3: User revokes donation consent.

6.5.1 Steps 1–3: Decide on revocation type

The user decides to revoke the donation consent (1) and is asked whether the revocation shall be temporary or permanent (2). Based on her answer (3), the process dove-tails.

6.5.2 Steps 4–12: Temporary revocation

Temporary revocation implies that we keep some information (namely the donor key pair) in order to restore the donor status later on, if the user wants to. We need the user's consent to retain that information. The corresponding consent steps 4–9 are essentially the same as the consent steps 2–7 in Figure 6.1. The different consent type χ vs. ξ shall illustrate that the consent text is different. Steps 10–12 introduce another consent type η which represents the user's actual intent to suspend or disable her donor capability. This concludes the steps required for temporary revocation.

6.5.3 Steps 10–18: Initialize permanent revocation

For permanent revocation we will delete the user's donor key pair. Hence, we do not need to ask for permission to keep it and can skip steps 4–9. Steps 10–12 are identical to the temporary revocation process and record the user's intent to disable her donor capability. Note, this covers only the formal intent, but proof (of the deletion of the donor key mapping at the donation service) is required before the consent service can accept a fresh future donor registration. Hence, the goal is to tell the donation service to drop the records containing the user's donation public key E^{Don} , and issue a cryptographic proof of this fact which can be checked by the consent service.

The donation key pair is retrieved from PHDP in steps 13 and 14. Step 15 creates a random nonce ρ which we will include in subsequent communication to prevent replay attacks. The revocation request we want to send in step 18 consists of the ciphertext of the donor public key E^{Don} , the nonce ρ , and the user's signature of the first two (16 and 17).

6.5.4 Steps 19–25: Delete donor key mapping

The donation service now decrypts the revocation request M (step 19) and verifies the included signature σ_ρ which proves that the request originated from the user (step 20). Next, all data associated with E^{Don} (for example, the ALP ID) is deleted in step 21. The deletion proof Π consists of the deletion timestamp τ , the nonce ρ , and the donation service's signature σ_Π of the two (steps 22 and 23). This proof is stored by the donation service (under the donor public key E^{Don}) for seven days. In case the communication with the client is disrupted at this point in time, it may never carry out the next steps 26–29. This would prevent the user from ever becoming a donor again (because the user never showed proof of the donor key mapping deletion

to the consent service). The caching of the proof Π and its signature enables the client to come back later and receive the proof.

6.5.5 Steps 26 – 29: Prove mapping deletion

The user now shows the deletion proof and its signature to the consent service (26). It will verify the signature (27), check that timestamp τ has not expired, and record that the user may apply as a donor again in the future (28). The proof is stored as invalidated until its timestamp expires to prevent replay attacks.

Finally, the client can delete the donor key pair (steps 30 and 31) and conclude the revocation process (32).

Part II

Formal treatment

7 Cryptographic notation

This chapter introduces short-hand notations that we use to accurately, yet concisely, define cryptographic protocols. We already implicitly used some notation in Chapter 2 and introduced some more in Section 3.1.

The main objectives for the formal notation are short names for cryptographic symbols (such as keys, messages or tags) and the encryption itself. Further, the notation should work at a whiteboard, on paper and have flavor that can be used in plain ASCII (for example in source code, Markdown or Slack). We use the following syntax to denote a cryptographic symbol or data element:

$${}^{(u)}s_b^A$$

It is a main symbol (here: s) which can be decorated with one index and up to two superscripts.

Symbol s

The main symbol is a single character that represents the type of object as listed in Table 7.1. A lowercase symbol indicates that the object is protection-worthy and must not be disclosed to the public. An uppercase symbol denotes an object which may (or must) be made public for certain protocols to work.¹

Superscript A

The superscript is used to indicate a subtype when used for certain symbols, like keys or passwords. See Table 7.2 for details.

Index b

We use the subscript for indexing or enumeration. The subscript itself can be complex or comma-separated or both. See the examples in Table 7.3.

¹This was inspired by the Go programming language where lowercase symbols are private while uppercase symbols get exported.

Symbol	Meaning
k	symmetric key
d	private key
E	public key
r	record
a	attachment
t	tag or token (clear from context)
S	salt value
w	password
m, v	other confidential message or value
M, V	other non-confidential message or value

Table 7.1: Different types of cryptographic symbols.

Symbol	Meaning
U	user-related
A	application-related
C	common key
D	data key
Don	donation key
N	attachment key
T	tag-related
F	function-derived (e. g. via a key derivation function)
S	service-related

Table 7.2: Superscripts indicating subtypes of cryptographic symbols.

Superscript (u)

We occasionally use this second superscript for tagging the symbol with a user ID u , if different users shall be differentiated explicitly. The superscript is written in parentheses.

Encryption

Typically, when denoting the encryption of some data x using a key k , some notation of the following is used:

$$\text{Enc}(k, x) \text{ or } \text{Enc}_k(x) \text{ or } E_k(x)$$

This notation can take up quite some screen real estate if many ciphertexts are part of the game (and this is the case for our platform). That is why we denote the ciphertext of x using key k as follows:

$$\left[x \right]_k$$

The precise cipher being used (symmetric or asymmetric) will be clear given the key notation.

ASCII flavor

We want to use the notation also for source code comments as well as in Slack communication. For this purpose we use the following ASCII rendering:

$$^{(u)}s_b^A \Leftrightarrow (u)sA_b$$

The main instances of the notation are listed in the following Table 7.3.

Symbol	ASCII	Meaning
d^U	dU	User private key (of implicit user)
$(u_1)d^U$	(u1)dU or (u_1)dU	User private key of user u_1
E^U	EU	User public key (of implicit user)
k_0^C	kC_0	Initial common key (of implicit user)
k^T	kT	Tag encryption key (of implicit user)
(d_1^A, E_1^A)	(dA_1, EA_1)	Application key pair of app ID 1
(d_s^A, E_s^A)	(dA_s, EA_s)	Application key pair for web app session s
$(d_{\text{Don}}^S, E_{\text{Don}}^S)$	(dS_Don, ES_Don)	Donation service key pair
w_P	w_P	User-chosen password
w_R	w_R	Recovery password
k_P^F	kF_P	Key derived from user password w_P
k_R^F	kF_R	Key derived from recovery password w_R

Symbol	ASCII	Meaning
$[d^U]_{k_P^F}$	[dU]kF_P	User private key encrypted using password-derived key
r_i	r_i	Record i
k_i^D	kD_i	Data key of record r_i
k_i^N	kN_i	Attachment key of record r_i
$[r_i]_{k_i^D}$	[r_i]kD_i	Ciphertext of record r_i encrypted using data key k_i^D
$[k_i^D]_{k_0^C}$	[kD_i]kC_0	Ciphertext of data key k_i^D encrypted using initial common key k_0^C
$a_{i,j}$	a_i,j or a_(i,j)	j -th attachment of record r_i
$[a_{i,3}]_{k_i^N}$	[a_i,3]kN_i	Ciphertext of 3rd attachment of record r_i
$t_{i,j}$	t_i,j or t_(i,j)	j -th tag of record r_i
$[t_{i,j}]_{k^T}$	[t_i,j]kT	Ciphertext of j -th tag $t_{i,j}$ (of record r_i) using tag encryption key k^T

Table 7.3: Examples of the cryptographic notation.

8 Account creation

When a new user registers an account with Data4Life, the following steps are carried out (see also Section 3.2 for a shorter but less precise description).

Client side

In the client application, the user enters the values shown in Table 8.1. A hash of the entered password is checked against the *haveibeenpwned* API [1] and in case of a positive result the user must choose another password. The same check is carried out when users want to change their password. The client application then computes and generates the data listed in Table 8.2. The data that gets sent to the server is listed in Table 8.3.

Server side

Most of the user registration payload becomes a row in the `users` table. The password hashes H_P and H_R are hashed again using *bcrypt* [18] to form H_P^* and H_R^* , respectively. The user public key and the encrypted initial common key become the zero permission

$$\pi_\emptyset := \left((u)E^U, \left\{ \left[(u)k_0^C \right]_{(u)E^U} \right\} \right).$$

See Figure 8.1 for a graph depicting the data dependencies of the user registration payload.

Material	Description
$(u)V_{\text{email}}$	User u enters e-mail address
$(u)_{WP}$	User u enters password

Table 8.1: Data required to be entered by the user when registering a new account.

Material	Description
$(^u)S_P$	User password salt
$(^u)k_P^E$	Derive key $(^u)k_P^E := \text{PBKDF2}((^u)w_P, (^u)S_P, N)$ with salt $(^u)S_P$ and N iterations
$(^u)w_R$	Generate recovery password (BIP-39 mnemonic)
$(^u)S_R$	Recovery password salt
$(^u)k_R^E$	Derive key $(^u)k_R^E := \text{PBKDF2}((^u)w_R, (^u)S_R, N)$ with salt $(^u)S_R$ and N iterations
$((^u)d^U, (^u)e^U)$	Generate user key pair
$(^u)k_0^C$	Generate initial common key
$(^u)k^T$	Generate tag encryption key

Table 8.2: Automatically generated data by the client.

Material	Description
$(^u)V_{\text{email}}$	User e-mail address
H_P	User password hash $H_P := \text{PBKDF2}((^u)w_P, 0, N)$
H_R	Recovery password hash $H_R := \text{PBKDF2}((^u)w_R, 0, N)$
$(^u)e^U$	User public key
$\left[(^u)d^U \right]_{(^u)k_P^E}$	User-password-encrypted user private key
$\left[(^u)d^U \right]_{(^u)k_R^E}$	Recovery-password-encrypted user private key
$(^u)S_P$	User password salt
$(^u)S_R$	Recovery password salt
$\left[(^u)k_0^C \right]_{(^u)e^U}$	Encrypted common key
$\left[(^u)k^T \right]_{(^u)k_0^C}$	Encrypted tag encryption key

Table 8.3: User registration payload that gets sent to the server.

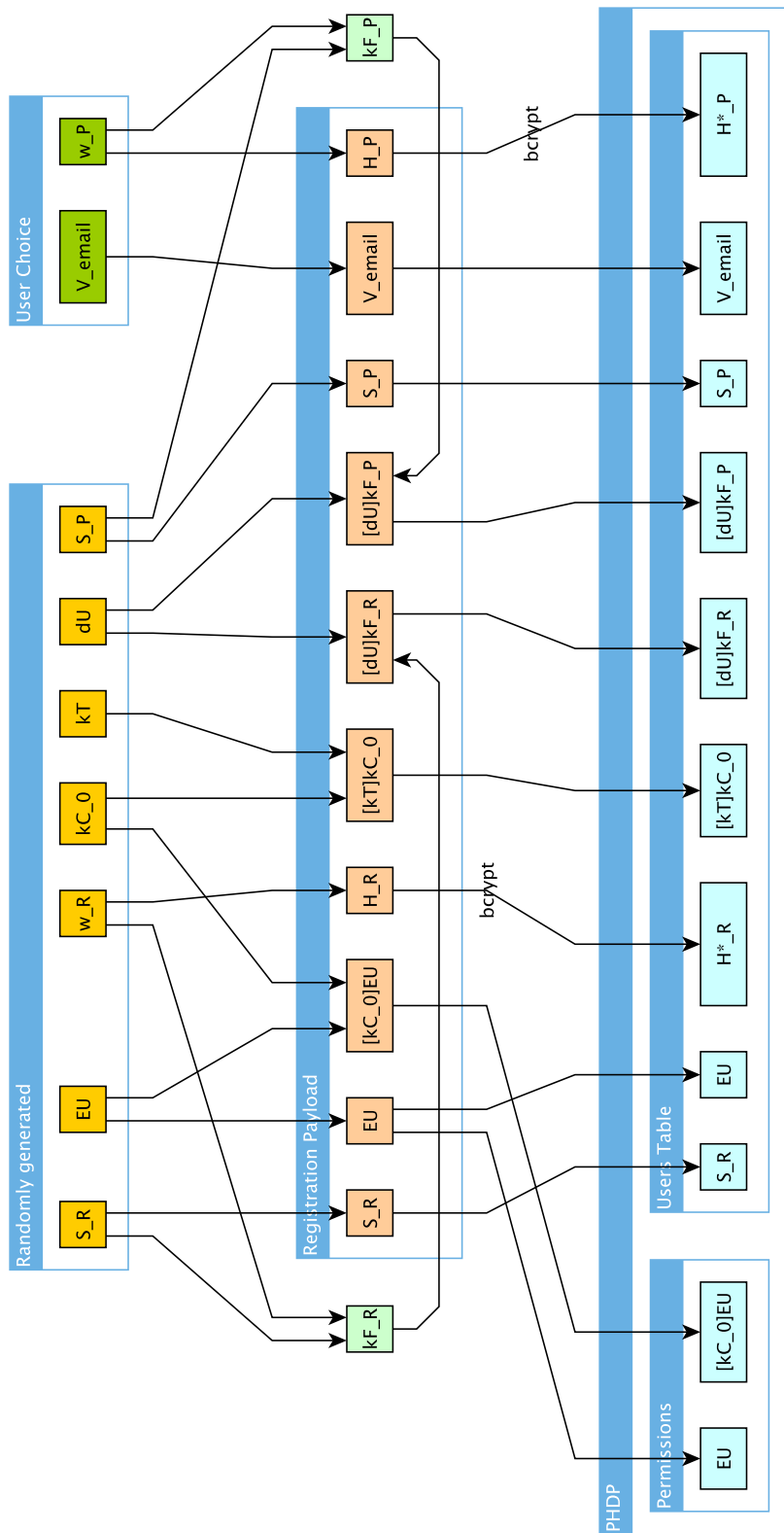


Figure 8.1: Dependencies of data elements needed for user account creation.

Part III

Appendices

A Bibliography

- [1] ‘;-have i been pwned? URL: <https://haveibeenpwned.com/> (visited on 07/30/2020) (cit. on pp. 17, 46).
- [2] Jean-Philippe Aumasson. *Serious Cryptography: A Practical Introduction to Modern Encryption*. USA: No Starch Press, 2017. ISBN: 978-1-59327-826-7 (cit. on p. 7).
- [3] *Authenticated encryption*. Wikipedia. URL: https://en.wikipedia.org/wiki/Authenticated_encryption (visited on 05/05/2020) (cit. on p. 11).
- [4] Mihir Bellare and Phillip Rogaway. “Optimal Asymmetric Encryption – How to Encrypt with RSA”. In: Springer-Verlag, 1995, pp. 92–111 (cit. on pp. 12, 52).
- [5] T. Dierks and E. Rescorla. *The Transport Layer Security (TLS) Protocol Version 1.2*. RFC 5246. RFC Editor, Aug. 2008. URL: <http://www.rfc-editor.org/rfc/rfc5246.txt> (cit. on p. 14).
- [6] *Documentation FHIR v4.0.1*. HL7 FHIR Website. URL: <https://www.hl7.org/fhir/documentation.html> (visited on 05/13/2020) (cit. on p. 23).
- [7] Morris Dworkin. *Recommendation for Block Cipher Modes of Operation: Galois/Counter Mode (GCM) and GMAC*. Tech. rep. SP 800-38D. NIST, 2007. URL: <https://doi.org/10.6028/NIST.SP.800-38D> (cit. on p. 10).
- [8] Morris J. Dworkin et al. *Advanced Encryption Standard (AES)*. Tech. rep. NIST FIPS 197. NIST, 2001. URL: <https://nvlpubs.nist.gov/nistpubs/FIPS/NIST.FIPS.197.pdf> (cit. on pp. 7, 52).
- [9] T. Elgamal. “A public key cryptosystem and a signature scheme based on discrete logarithms”. In: *IEEE Transactions on Information Theory* 31.4 (1985), pp. 469–472 (cit. on p. 12).
- [10] Niels Ferguson, Bruce Schneier, and Tadayoshi Kohno. *Cryptography Engineering: Design Principles and Practical Applications*. Wiley Publishing, Inc., 2010. ISBN: 978-0-470-47424-2 (cit. on pp. 5, 7, 9, 53).
- [11] D. Hardt. *The OAuth 2.0 Authorization Framework*. RFC 6749. RFC Editor, Oct. 2012. URL: <http://www.rfc-editor.org/rfc/rfc6749.txt> (cit. on pp. 20, 28).

- [12] M. Jones, J. Bradley, and N. Sakimura. *JSON Web Token (JWT)*. RFC 7519. RFC Editor, May 2015. URL: <http://www.rfc-editor.org/rfc/rfc7519.txt> (cit. on p. 19).
- [13] B. Kaliski. *PKCS #5: Password-Based Cryptography Specification Version 2.0*. RFC 2898. RFC Editor, Sept. 2000. URL: <http://www.rfc-editor.org/rfc/rfc2898.txt> (cit. on p. 19).
- [14] Alfred J. Menezes, Paul C. Van Oorschot, and Scott A. Vanstone. *Handbook of Applied Cryptography*. 1st. USA: CRC Press, Inc., 1997. ISBN: 978-0-84-938523-0 (cit. on pp. 7, 12).
- [15] *Mnemonic code for generating deterministic keys*. Github.com. URL: <https://github.com/bitcoin/bips/blob/master/bip-0039.mediawiki> (visited on 05/07/2020) (cit. on p. 17).
- [16] National Institute of Standards and Technology. *Data Encryption Standard (DES)*. Tech. rep. FIPS PUB 46-2. NIST, 1993. URL: <https://doi.org/10.6028/NIST.FIPS.46-2> (cit. on pp. 7, 52).
- [17] Christof Paar and Jan Pelzl. *Understanding Cryptography: A Textbook for Students and Practitioners*. 1st. Springer Publishing Company, Inc., 2009. ISBN: 978-3-642-04100-6 (cit. on p. 7).
- [18] Niels Provos and David Mazières. “A Future-Adaptive Password Scheme”. In: *Proceedings of the Annual Conference on USENIX Annual Technical Conference*. ATEC '99. Monterey, California: USENIX Association, 1999, p. 32 (cit. on pp. 19, 46).
- [19] R. L. Rivest, A. Shamir, and L. Adleman. “A Method for Obtaining Digital Signatures and Public-Key Cryptosystems”. In: *Communications of the ACM* 21.2 (Feb. 1978), pp. 120–126. ISSN: 0001-0782. DOI: [10.1145/359340.359342](https://doi.org/10.1145/359340.359342). URL: <https://doi.org/10.1145/359340.359342> (cit. on p. 12).
- [20] Bruce Schneier and Phil Sutherland. *Applied Cryptography: Protocols, Algorithms, and Source Code in C*. 2nd. USA: John Wiley & Sons, Inc., 1995. ISBN: 978-0-471-12845-8 (cit. on p. 7).

B Glossary

Term	Abbrev.	Definition
Advanced Encryption Standard	AES	Symmetric cipher for encryption and decryption of data [8]
Analytics Platform	ALP	Data analysis platform for de-identified healthcare data
authenticated encryption		A cipher which does not only provide confidentiality, but also message authenticity.
authenticated encryption with associated data	AEAD	Authenticated encryption where some additional data can be added which is not encrypted, but covered by the authentication tag.
authentication tag		Cryptographic checksum that guarantees the integrity of a ciphertext.
common key		Symmetric key used to encrypt data keys and attachment keys.
Data Encryption Standard	DES	Outdated symmetric cipher for encryption and decryption of data [16]
end-to-end encryption	E2EE	Encryption and decryption of healthcare data occur at the client. The server never gets in contact with unencrypted (plaintext) healthcare user data.
initialization vector	IV	Random data that is used as the “zeroth” plaintext block to start encryption.
mode of operation		Algorithm to iteratively apply a (fixed-size) block cipher to arbitrary input message lengths.
Optimal asymmetric encryption padding using RSA	RSA-OAEP	Padding scheme used in conjunction with RSA that is proved secure against certain attack types [4].

Term	Abbrev.	Definition
padding		Process of extending a plaintext to give it a length that is a multiple of the supported block size of the used symmetric cipher. Also, it is used with asymmetric encryption to randomize the plaintext in order to avoid certain attack scenarios [10].
Personal Health Data Platform	PHDP	Data storage platform using end-to-end encryption for secure healthcare data
tag encryption key		Symmetric key used to encrypt records tags.
user private key		Private key of an RSA key pair created at user registration and valid throughout the account lifetime.
user public key		Public key of an RSA key pair created at user registration and valid throughout the account lifetime.

C Change history

Version	Date	Remarks
1.0.0	2020-05-13	Initial revision
1.1.0	2020-05-20	Add data sharing, crypto notation and examples
1.1.1	2020-05-25	Adjust wording, fix typo
1.1.2	2020-05-28	Add acknowledgments
1.1.3	2020-06-05	Fix typos
1.1.4	2020-06-06	Adjusted profession wording
1.2.0	2020-06-08	Streamline chapters
1.2.1	2020-06-11	Incorporate feedback
1.2.2	2020-07-08	Fix typo, fix reference
1.3.0	2020-07-09	Add data donation chapter
1.3.1	2020-07-30	Add information on password checks
1.3.2	2020-08-17	Fix wording in data ingestion chapter

D Acknowledgments

This work has received funding from the European Union's Horizon 2020 research and innovation program under the grant agreement No. 826117: Smart4Health – Building a citizen-centered EU-EHR exchange for personalized health.



DATA4LIFE.CARE

Copyright 2020 D4L data4life gGmbH. All rights reserved.

data
4life