



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 826078.

Privacy preserving federated machine learning and blockchaining for reduced cyber risks in a world of distributed healthcare



Deliverable D6.5
“Mechanisms for removing sensitive information from the blockchain”

Work Package WP6
“Blockchains and user right management”

Disclaimer

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 826078. Any dissemination of results reflects only the author’s view and the European Commission is not responsible for any use that may be made of the information it contains.

Copyright message

© FeatureCloud Consortium, 2022

This deliverable contains original unpublished work except where clearly indicated otherwise. Acknowledgement of previously published material and of the work of others has been made through appropriate citation, quotation or both. Reproduction is authorised provided the source is acknowledged.

Document information

Grant Agreement Number: 826078		Acronym: FeatureCloud	
Full title	Privacy preserving federated machine learning and blockchaining for reduced cyber risks in a world of distributed healthcare		
Topic	Toolkit for assessing and reducing cyber risks in hospitals and care centres to protect privacy/data/infrastructures		
Funding scheme	RIA - Research and Innovation action		
Start Date	1 January 2019	Duration	60 months
Project URL	https://featurecloud.eu/		
EU Project Officer	Christos MARAMIS, Health and Digital Executive Agency (HaDEA) - Established by the European Commission, Unit HaDEA.A.3 – Health Research		
Project Coordinator	Jan BAUMBACH, UNIVERSITY OF HAMBURG (UHAM)		
Deliverable	D6.5 “Mechanisms for removing sensitive information from the blockchain”		
Work Package	WP6 “Blockchains and user right management”		
Date of Delivery	Contractual	31/12/2022 (M48)	Actual 19/12/2022 (M48)
Nature	Report	Dissemination Level	Public
Lead Beneficiary	05 SBA		
Responsible Author(s)	Walid Fdhila, SBA		
Keywords	Blockchain, GDPR, Federated machine learning, auditability		

Table of Content

1	Objectives of the deliverable based on the Description of Action (DoA)	4
2	Executive Summary	4
3	Introduction (Challenge)	4
4	Methodology	5
5	GDPR and Blockchain	5
5.1	Data Subject, Controller and Processor	5
5.2	Personal data and anonymization	6
5.3	GDPR in the Context of Blockchain Data and Identifiers	8
5.4	Data processing	10
5.5	Data Controllers and processors in Blockchain	10
5.6	Rights to Rectification and Erasure	10
6	The FeatureCloud Blockchain and the Right to Erasure	11
6.1	Overall FeatureCloud Architecture	11
6.2	FeatureCloud Blockchain and GDPR compliance	14
6.2.1	Personal Data and Data processing in FeatureCloud	14
6.2.2	Data Controllers, and Processors In FeatureCloud Blockchain	16
6.2.3	How does FeatureCloud Blockchain meet the right to rectification and erasure	17
7	Conclusion	21
8	References	22
9	Table of acronyms and definitions	22

1 Objectives of the deliverable based on the Description of Action (DoA)

The objective of this deliverable is based on the description of action, which incorporates mostly aspects from Objective 3 “To develop mechanisms for selective deletion of sensitive information from the blockchain-mechanism, thus allowing practical consent revocation as demanded by the GDPR (Task 3)”. The corresponding task in this work package is Task 3 (SBA, RI): “Current blockchain mechanisms are designed for not allowing alterations, i.e., protecting the integrity of past transactions, either by final or by probabilistic consent mechanisms. While this is a vital feature, it is important to unlink the sensitive information from the direct blockchain internal mechanisms, i.e., allowing for selective deletion of sensitive information. SBA will thus research into blockchain mechanisms that possess consensus mechanisms that allow for the selective deletion of information, in case all partners support the decision. The developed concepts will be tested for actual performance by SBA.”

2 Executive Summary

Deliverable D6.5 builds on the legal requirements for consent management collected in D6.3, which have been heavily discussed with the legal experts within the FeatureCloud consortium and upon which a GDPR-friendly solution design was proposed in accordance with Article 25 GDPR on data protection by design and by default. The solution also improves trust of the auditability and accountability processes of the data controllers and processors as defined in Article 30, GDPR. In D6.5, we will first discuss how GDPR applies in the case of Blockchain technology with a special focus on the rights to data rectification and erasure and then demonstrate how the solution design proposed for FeatureCloud complies with the regulation.

3 Introduction (Challenge)

Over the last decade, Blockchain and Distributed Ledger technologies (DLT) have emerged as a key enabler for digital transformation by means of innovative ways of data recording and sharing without central control. Because of their reliance on distributed systems, sophisticated consensus mechanisms, and advanced cryptographic techniques, DLTs offer a range of nice properties such as resilience, immutability, non-repudiation and integrity. However, while these very characteristics are required to ensure security, auditability and transparency, they can be problematic in regard to privacy, confidentiality and data protection requirements such as enabling the removal of undesirable content or otherwise deleting or changing the recorded transaction history. This tension between the immutability of DLTs and the rights for rectification and erasure as stipulated by the European Union' General Data Protection Regulation (GDPR) in articles 16 and 17 raises concerns as to which extent GDPR will affect blockchain-based use cases or blockchain adoption. While many would consider blockchain and GDPR as fundamentally incompatible, others believe that GDPR-compliant blockchain solutions that respect the fundamental principles of data protection are possible if designed and implemented properly. Others also regard blockchain technology as means to achieve some of the GDPR objectives such as transparency and data control through well-defined and enforced access rights (e.g., using smart contracts). Furthermore, the existence of different forms of DLTs that employ various technical designs and governance set-ups requires GDPR to be assessed on a use-case basis rather than generalised to all blockchains.

The FeatureCloud project aims at providing privacy-by-design solutions for research on medical data using federated machine learning of healthcare data. Instead of collecting healthcare data in a

central repository, machine learning models are trained locally at each data provider site, and the output is collected and aggregated using secure multi-party computation or other technologies. This avoids legal and privacy issues encountered in centralised learning as the data remains within the institution and jurisdiction where it was collected. In previous Deliverables D6.1 to D6.4, we have demonstrated that blockchain technology (BT) can indeed be an essential tool to improve the auditability of federated machine learning processes, and managing patient consents and access controls. Thereby, an immutable audit trail is created that can be used to detect deviations in retrospect. The solution we proposed is based on technical but also legal requirements from GDPR (cf. Deliverable D6.3), and is privacy-and-compliant-by-design. In the context of this deliverable D6.5, we will discuss how our design complies with data protection regulations and more specifically with user rights to erasure of personal data. First, we will present some of the key concepts and challenges from GDPR, then we will outline how they are addressed in the current design.

4 Methodology

Articles 16 and 17 of the GDPR define the circumstances under which an individual has the right to erase or rectify personal data. While it would appear that the individual ability to exercise those rights would contradict with the blockchain properties, with careful design considerations, the advantages of BT can be leveraged while avoiding these privacy issues. As such, FeatureCloud follows a privacy-by-design methodology by proactively embedding privacy and data protection into the solution design. The legal and technical requirements collected and analysed in D6.3, allowed us to provide a blockchain-based audit trail with data minimization in mind, where only the minimum amount of data (commitments) necessary for audit and consent management is stored on-chain. Deliverable D6.5 focuses on the FeatureCloud blockchain solution, and demonstrates how the latter complies with the right to erasure. First, we conducted an analysis of what constitutes personal data under the GDPR and explored current interpretations of the GDPR in the context of blockchain. Then, we focused on personal data in the FeatureCloud blockchain and demonstrated how patient re-identification based on data committed to the blockchain is reasonably unlikely. By keeping the actual data off-chain and employing appropriate and secure cryptographic techniques, it becomes possible to comply with the GDPR requirements. In particular, the approach relies on the irreversibility of the commitments or the reconstruction of the original data from the corresponding hashes due to the large domain space for data representations.

5 GDPR and Blockchain

The GDPR came into effect on May 25th, 2018 to protect the rights of natural persons and regulate how their personal information is collected, shared and used. In the following, we define some of the key terms that will be used along the deliverable and the legal description of specific rights with which our solution design needs must comply. We will also discuss how they are interpreted in the case of DLTs. More specifically, we will focus on the rights to rectification and erasure.

5.1 Data Subject, Controller and Processor

- According to article 4(1) of the GDPR, a **data subject** is an identified or identifiable natural person, whose personal data is processed, where processing means any operation or set of operations which is performed on personal data or on sets of personal data, whether or not by automated means, including recording, mere storage and also destruction of that data (see in detail chapter 4.4. below).

- The GDPR definition of a **controller** is “the natural or legal person, public authority, agency or other body which, alone or jointly with others, determines the purposes and means of the processing of personal data” (see Article 4(7) GDPR).
Note that it is possible to have more than one controller that jointly determine the purpose and means of the processing, namely joint controllers (see Article 26 GDPR).
- Finally, in article 4(8) of the GDPR, a **processor** means a natural or legal person, public authority, agency or other body which processes personal data on behalf of the controller”.

5.2 Personal data and anonymization

The scope of GDPR application is fundamentally defined by what qualifies as personal data. Determining whether or not a piece of data is personal is essential but also not trivial given the broadness of GDPR definitions and the various factors that enter into play, mainly driven by new technological developments, e.g., new anonymization techniques, new cryptographic algorithms that anonymize or pseudo-anonymize personal data. For example, it is not always easy to determine to which extent a given anonymization technique is considered to be GDPR compliant and the anonymized data as not personally identifying, e.g., substituting a name for an identifier may still be representative of that person and thus intrusive to personal privacy¹.

Article 4(1) GDPR, for example, qualifies as **personal data** “any information relating to an identified or identifiable natural person such as a name, an identification number, location data, an online identifier or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person”. In relation to healthcare in general, and FeatureCloud in particular, Article 4(15) of the GDPR defines **data concerning health** as personal data related to the physical or mental health of a natural person, including the provision of health care services, which reveal information about his or her health status, which should, put more precisely by Recital 35 of the GDPR, include “all data pertaining to the health status of a data subject which reveal information relating to the past, current or future physical or mental health status of the data subject”. According to Article 9(1) GDPR data concerning health belongs to the special categories of personal data which may only be processed in one of the particularly restrictive cases set out in Article 9(2) GDPR.

What constitutes anonymous Data?²

GDPR distinguishes between anonymization and pseudonymization, and only classifies anonymized data as non-personal data, and therefore out of GDPR scope. However, the uncertainty about what constitutes anonymous data continues to be increasingly burdensome and subject to debate. Article 4(5) GDPR defines pseudonymization as “the processing of personal data in such a manner that the personal data can no longer be attributed to a specific data subject without the use of additional information, provided that such additional information is kept separately and is subject to technical and organisational measures to ensure that the personal data are not attributed to an identified or identifiable natural person”. This underlines that pseudonymized data remains personal data, but does not dismiss pseudonymization as means to secure data as stated in Article 5(f) GDPR.

Anonymized data, on the other hand, cannot be associated with specific individuals. In contrast to pseudonymization, anonymized data is irreversible and the likelihood of an individual’s re-identification through the use of data analysis or similar techniques is negligible in practice. This means that it is no longer possible to reconstruct the original personal data from the anonymized data, thereby falling out of the GDPR scope. Recital 26 GDPR refers to anonymous information as being “information which does not relate to an identified or identifiable natural person” and anonymised information as “personal data rendered anonymous in such a manner that the data

¹ Edwards L (2018), Law, Policy and the Internet, Oxford: Hart Publishing, 85

² Seen in more detail already D2.1 “Risk Assessment Methodology”, Chapter 5 “Privacy Law Analysis”.

subject is not or no longer identifiable”. Recital 26 GDPR defines a risk-based method to determine whether or not a data is personal (see figure 1). However, the difficult determination of what constitutes a ‘reasonable likelihood’ of identification further burdens practitioners’ work (Finck and Pallas, 2020).

Identifiability

In the current formulation of the GDPR, it remains unclear from whose perspective the likelihood of identifiability ought to be assessed, i.e., the data controller or any third party capable of identifying the data subject.³ For instance, the Article 29 Working Party⁴, and in line with the risk-based approach of GDPR, outlined three main criteria that help determine whether de-identification occurs as follows:

- **singling out:** it is still possible to single out an individual
- **linkability:** it is still possible to link records relating to an individual, and
- **inference:** information concerning an individual can still be inferred.

In the context of blockchain, and FeatureCloud in particular, it is very important to have a clear understanding of which data that is frequently used in relation to blockchain qualifies as personal data (e.g., public keys, hashes) and whether pseudonymisation measures can produce anonymous data. Indeed, pseudonymous data on a blockchain can, in principle, be related to an identified or identifiable natural personal through singling out, inference or linkability (Finck, 2019a).

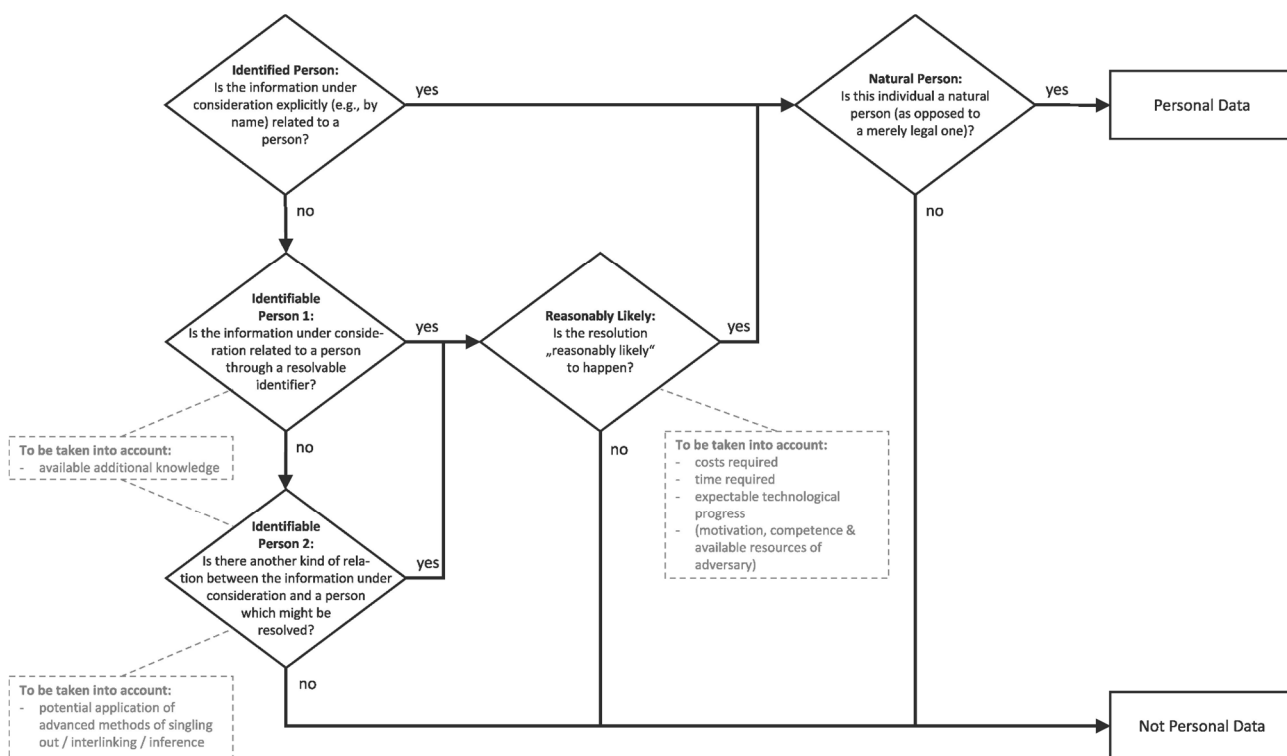


Figure 1. Diagram for assessing data to person relatedness under the GDPR (From (Finck and Pallas, 2020))

³ Cf. on this question the pre-GDPR Breyer decision (ECJ 19.10.2016, C-582/14).

⁴ A29WP Opinion on Anonymisation Techniques (2014)

5.3 GDPR in the Context of Blockchain Data and Identifiers

(The following concepts are important in determining whether leaving hashes on blockchain while deleting the off-chain data still complies with the right to erasure). Also, In the context of consent management, it would be crucial to determine where the use of public keys to update consents would be a correlating factor (although we use local identifiers per participant).

Cryptographic Hashes as Personal Data

A hash function (H) is a mathematical *deterministic* function that takes a message (m) of arbitrary but finite size and outputs a fixed size hash (a.k.a. digest). This means that the same input will always yield the same output. The properties of cryptographic hash functions allow us to uniquely identify data based on the computed hash value with high probability. Hereby, this probability relates to the specific security parameters of the hash function, in particular its collision resistance and second pre-image resistance. Furthermore, the property of pre-image resistance allows cryptographic commitments where a party commits to a secret message by only revealing its hash. At the time the secret is revealed anyone can compute the hash function with the secret as input and verify if the computed hash matches the original commitment. To ensure that pre-image resistance is achieved, the input message must be of sufficient entropy, i.e., the input space of possible messages must be large enough to render any guessing attacks infeasible. In practice this means that data with low entropy, e.g., dates of birth, must be amended, usually by concatenating to the message some random secret, to prevent an adversary from being able to guess the pre-image through brute force.

In the context of blockchain, hash functions are one of the most extensively used cryptographic schemes. Besides their main role in ensuring blockchain data integrity and their use in mining (specifically for proof of work PoW), hash functions are also employed as subject identifiers or commitments of off-chain data.

With respect to GDPR, hash functions have been used as pseudonymization techniques to remove the explicit links to data subjects from data (as an **identifier replacement**), as to prevent content-based re-identification (**content replacement**). For example, hashing a patient named “Bob” (i.e., ID replacement) through the hashing function SHA256 will always yield the result “81b637d8fcd2c6da6359e6963113a1170de795e4b725b84d1e0b4cf9ec58ce9”. However, despite its pre-image resistance (i.e., deriving the name from the hash is not possible), the hash output is not necessarily anonymous for many reasons. Indeed, anyone who knows that Bob data may be in a specific dataset could run the same algorithm on “Bob” and easily identify all transactions Bob was involved in. Moreover, it is also possible to correlate data from multiple datasets that belong to the same identifier, thereby violating the assessment test as defined in Recital 26 GDPR. This is particularly possible and easy to achieve if the input space is relatively small and in the absence of additional measures that render brute forcing infeasible (e.g., hash salting (Gauravaram, 2012)).

In fact, assessing whether brute forcing hash-based re-identification is “reasonably likely” depends on the time frame, the number of attempts and cost required to re-identify a subject. As such, using hashes of patient identifiers (e.g., hashes of social security numbers or names) as identification mechanism for data subjects may appear weak and insufficient for protecting them against privacy intrusions. For example, a graphic card such as the NVIDIA GeForce RTX 4090 has a hash rate power of 5.2 Gh/s, which means it can perform 5.5 billion hash guesses per second, and only costs under 2k euros. Top bitcoin miners have hash rates in the order of EH/s (quintillions of hashes per second). Therefore, subject re-identification can be “reasonably likely” as soon as the input space is relatively small.

Similar to identifiers, content re-identification refers to the ability of identifying a data subject on the basis of the content data (e.g., address, genetics). Unlike identifiers, and unless the content data is small (e.g., just addresses, or birthdates), content data is unlikely to be re-identified from the corresponding hashes. In particular, if the input data set is large and additional techniques are employed such as aggregation, k-anonymity and differential privacy.

For both hash-based identifiers and content, common practices consist in adding random nonces to the input data before hashing. This increases the input domain space and decreases the likelihood of the data being brute forced, e.g., peppering and salting (Gauravaram, 2012). Note that peppering offers stronger privacy guarantees as it relies on a secret key (the pepper) as additional input. This said, assessing whether a re-identification is “reasonably likely” can still be not straightforward, and should be done on a case-by-case basis.

Public Keys and GDPR

In public-key cryptography (also known as asymmetric cryptography), a pair of keys (namely public and private keys) are generated using a cryptographic scheme. Both keys are mathematically correlated in a way that only the private key decrypts what a corresponding public key has encrypted and vice versa. Together, they can be used to encrypt, decrypt and digitally sign messages. The advantage of asymmetric cryptography is twofold. First, one does not have to trust that other parties keep any shared secret key secure. Second, it allows to readily and non-interactively share the necessary information, i.e., the public key, for encrypted communication, or for verification that a digital signature was created from the private key corresponding to a particular public key.

Public keys are considered as pseudonymous identifiers and may qualify as personal data subject to GDPR (Finck and Pallas, 2020). In the bitcoin blockchain for example, it was further demonstrated that using sophisticated analysis techniques (e.g., graph analysis), it is possible to correlate a public key with the corresponding natural person (Ghesmati et al., 2022a, 2022b). There is indeed apparent similarity with IP addresses which also classify as personal data as the subject can indeed be identified by either the ISP or the services they interact with.⁵

Encrypted Data

Similar in principle to hashing, encryption is a process of encoding data, where an input data is transformed into a ciphertext. Encryption is however a two-way process, which means that it is also possible to decrypt the ciphertext using a secret key to reproduce the original data. The purpose of encrypting data is to maintain the confidentiality of sensitive data. The likelihood of decrypting the data without possessing the secret key heavily depends on the encryption scheme employed. In the past, many of what had been considered as strong encryption schemes at the time are now broken and no longer used (e.g., data encryption algorithm DES).

Article 32(1) GDPR emphasises encryption as an appropriate measure to secure personal data by controllers and processors and protect the confidentiality of subjects’ data. Given the fact that anyone with access to the encryption key (legally or through theft) may have the data decrypted, thereby identifying the subject, encrypted data would be considered as personal data. The EU Blockchain Observatory and Forum also approves this finding and concludes that even if strong encryption is employed on personal data, the result is surely pseudonymous and not anonymous, rendering it in the scope of GDPR⁶.

⁵ Cf. ECJ 19.10.2016, C-582/14, *Breyer*.

⁶ https://www.eublockchainforum.eu/sites/default/files/reports/20181016_report_gdpr.pdf

5.4 Data processing

Article 4(2) GDPR defines processing as “any operation or set of operations which is performed on personal data or on sets of personal data, whether or not by automated means, such as collection, recording, organisation, structuring, storage, adaptation or alteration, retrieval, consultation, use, disclosure by transmission, dissemination or otherwise making available, alignment or combination, restriction, erasure or destruction”.

As this is the widest possible definition (“any operation”), adding and storing personal data as well as any further processing on a distributed ledger may fall in the scope of GDPR and constitutes personal data processing (Finck, 2019a).

Do Smart contracts count as data processing?

In a nutshell, smart contracts are automated programs that reside and run on Blockchain. They are code, which embed a set of predetermined conditions that, when met, trigger outcomes. Their application goes beyond the realm of crypto currencies and includes supply chain management and insurance claim management, and they may process personal data.

Given that Article 22(1) GDPR restricts “solely automated individual decision-making and profiling”, then understanding whether or not smart contracts may fall into this scope becomes crucial. (Finck, 2019b) and (Finck, 2019a) conclude that smart contracts may indeed, in some circumstances, qualify as a form of solely automated data processing under Article 22(1) GDPR unless they meet the requirements of Article 22(2) and implement the safeguards of Article 22(3) GDPR. Article 22(2)(c) GDPR, for example, allows automated data processing where it is based on the data subject's explicit consent, thereby providing the possibility to lawfully operate smart contracts that process personal data. This, of course, has to be complemented with safeguarding measures such as the right to human intervention of Article 22(3) GDPR.

5.5 Data Controllers and processors in Blockchain

As described in Section 5.1, under GDPR data controllers and processors are the entities that are held responsible for storing, securing and processing personal data and complying with the data protection regulations. Article 5 GDPR (2) mandates that the controller shall be responsible for, and be able to demonstrate compliance with, paragraph 1 (‘accountability’). Therefore, and in order to understand how to apply GDPR on a case-by-case basis, it is first necessary to identify who acts as the data controller and who acts as the data processor.

According to (Finck, 2019), the GDPR's broad definition of controllership has far-reaching implications for personal data processing based on DLT. The various actors and roles involved in the different layers of blockchain, in addition to mining participation dynamicity in certain blockchain designs (such as in permissionless blockchains), makes it further difficult to determine who are the controllers and processors as many actors influence the determination of the means of processing. Moreover, whether it is a public permissionless blockchain or private permissioned blockchain, governance as well as the governing entities influence differently the modalities of processing. The CNIL, for example, outlines that smart contract developers may be considered as joint controllers if they actively participate in data processing.

5.6 Rights to Rectification and Erasure

There is a dichotomy between user privacy/data protection requirements and the construction of a blockchain as a transparent and verifiable immutable ledger of transactions. On the one hand, BT promises to offer compelling characteristics and properties that can be leveraged, e.g., tamper resistance, high reliability, openness, and distributed or even decentralised trust. For instance, BT

can be used to realise global data sharing and data traceability systems, where these advantages make it possible to build larger scale, higher quality, and auditable global decentralised data platforms. On the other hand, the aforementioned properties of BT also present fundamental challenges in respect to ensuring some user rights such as enabling the removal of undesirable content or otherwise deleting or changing the recorded transaction history. The ability to both withhold and even delete data due to privacy and regulatory requirements constitutes a necessity for compliant systems. At first, it would appear that in such cases incorporating BT is not an ideal approach. However, with careful design considerations the advantages of BT can be leveraged while avoiding these privacy issues.

The right to rectification

Article 16 of the GPR states that “the data subject shall have the right to obtain from the controller without undue delay the rectification of inaccurate personal data concerning him or her. Taking into account the purposes of the processing, the data subject shall have the right to have incomplete personal data completed, including by means of providing a supplementary statement”.

The right to erasure

Also known as the right to be forgotten, this gives individuals the right to ask controllers to delete their personal data following grounds such as unlawful processing of the data or the data is no longer necessary for the purpose for which they were collected. Therefore, Article 17 GDPR also outlines the circumstances under which an individual has the right to have their personal data erased. According to Article 17, “the data subject shall have the right to obtain from the controller the erasure of personal data concerning him or her without undue delay and the controller shall have the obligation to erase personal data without undue delay”.

Furthermore, “where the controller has made the personal data public and is obliged pursuant to paragraph 1 to erase the personal data, the controller, taking account of available technology and the cost of implementation, shall take reasonable steps, including technical measures, to inform controllers which are processing the personal data that the data subject has requested the erasure by such controllers of any links to, or copy or replication of, those personal data”.

It should be noted that the right to erasure does not apply to all circumstances and there are exceptions, in which an organisation’s right to process someone’s data may override the right to erasure such as “the data being processed is necessary for public health purposes and serves in the public interest” or “for archiving purposes in the public interest, scientific or historical research purposes or statistical purposes in accordance with Article 89(1) in so far as the right referred to in paragraph 1 is likely to render impossible or seriously impair the achievement of the objectives of that processing”.

6 The FeatureCloud Blockchain and the Right to Erasure

In the following, we will first briefly introduce the design solution and overall FeatureCloud architecture, then we will discuss how GDPR and more specifically the right to erasure is handled by the design.

6.1 Overall FeatureCloud Architecture

In FeatureCloud, patient data and consents are collected and stored locally in secure databases at each participant site (e.g., hospital). When a participant is invited to be part of a federated machine learning (ML) study of healthcare data, the ML algorithm is executed locally and only the output

model is provided (federated learning, FL). This means that patient data and consents never leave the site storage. As shown in Figure 2, a participant needs to run the so-called FeatureCloud controller, which manages the local execution. On the coordinator side, the controller also orchestrates execution and instructs the participants’ controllers to ensure a globally synchronous execution.

The app containers cannot communicate directly with each other to restrict Internet access for security reasons. Instead, the controllers pass through the traffic, which uses a separate relay server. The FeatureCloud system includes the global API where project details are stored, and the AI store where FeatureCloud apps can be fetched from. Apart from aggregated model parameters, nothing related to patients, their consent or, most importantly, their raw data, leaves a participant’s or coordinator’s site. All data required by the apps are selected by the participants beforehand.

In previous deliverables D6.1 to D6.4, we highlighted the potential of utilising blockchain technology (BT) in FeatureCloud where ensuring data privacy and compliance constitutes a necessity. As such, we proposed a system design that uses BT to secure FL processes and improve their auditability. A key challenge for the design is to ensure accountability of the involved actors and prove that only legitimate data was used in the FL process, thereby enabling detect wrong doings such as the use of non-consented data, or fake patient data, identities and consents. This will help improve the integrity of FL studies and tackle the data poisoning problem in the presence of Byzantine actors where data protection is an issue.

More specifically, in deliverable D6.4, we described a design that relies on a permissioned private blockchain based on both legal and technical requirements collected in D6.3. The solution is based on hyperledger fabric and two particular smart contracts were implemented and deployed to manage ML studies and consents. Access controls were defined such that only authorised entities have access to the ledger data (with different restricted views on data), and can invoke smart contract functions.

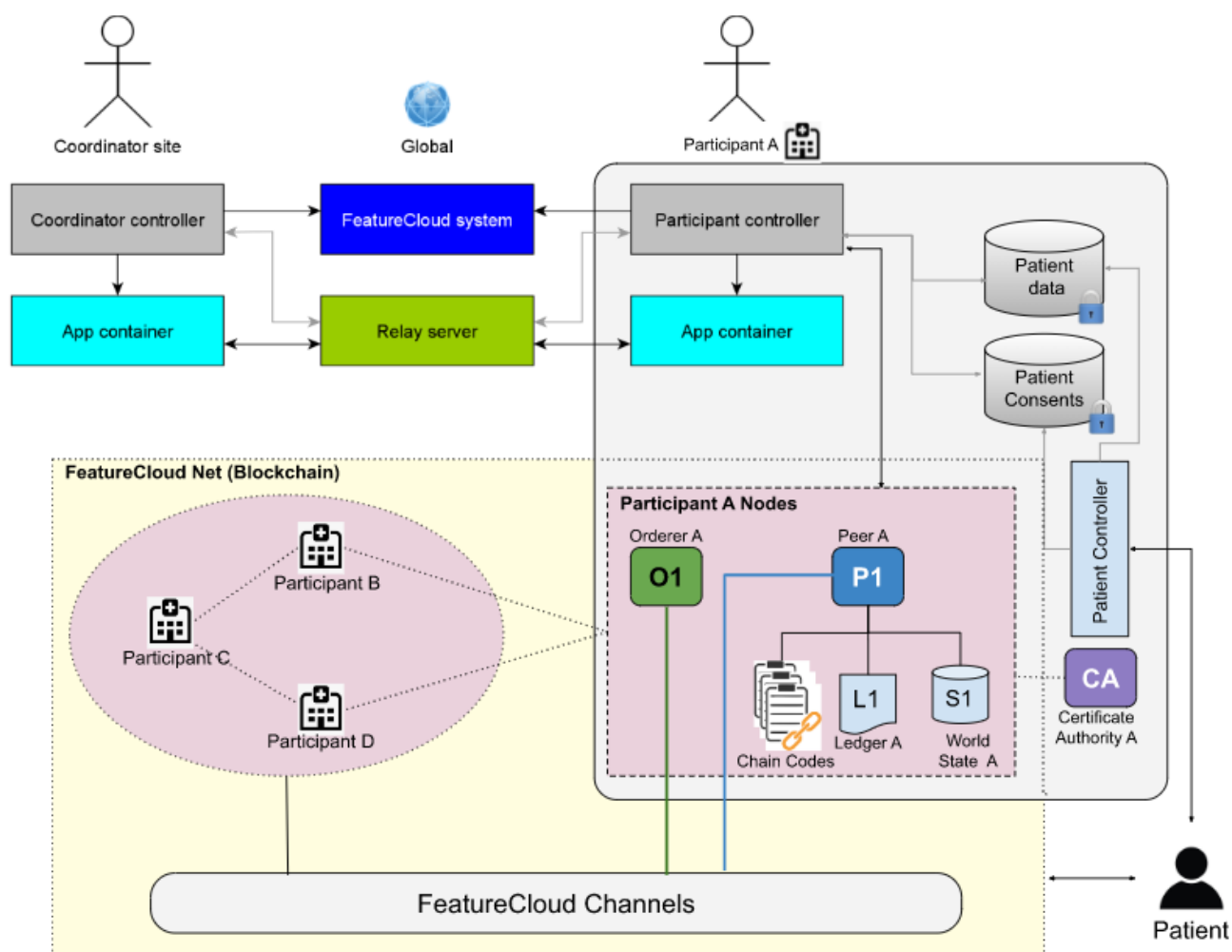


Figure 2. FeatureCloud Overall Architecture

Note that identities of all entities participating in FeatureCloud are legally binding. The FeatureCloud blockchain network is composed of different actors from the healthcare sector, which are onboarded following specific onboarding rules, and which constitute the trust framework (also referred to as governance framework). The latter is a network of trusted authorities TA (may initially be health ministries of EU countries) that defines the governance rules for all stakeholders and enables legally binding relationships (e.g., onboarding policies, applicable regulations, governance rules, protocols). Members of the trust framework are also responsible for onboarding new members to FeatureCloud (e.g., hospitals, pharmaceutical companies, test labs, insurances) with specific access rights (e.g., writing to the ledger, specific function invocations), and monitoring their compliance with the rules. All accredited members in this trust framework may act as both validator nodes and identity providers.

A participant (e.g., hospital) willing to participate in the blockchain infrastructure (running blockchain nodes), first needs to be accredited by the trust framework. However, this does not constitute a mandatory step to use the infrastructure. The participant may be given access to use the FeatureCloud blockchain as an audit trail for ML studies or for consent management without having to run its own nodes. As depicted in Figure 2, a participant in the FeatureCloud blockchain infrastructure needs to run one or several peer nodes and one or more orderer nodes. Peer nodes are responsible for endorsing transactions and hosting smart contracts (i.e., chain codes). Each peer node maintains a copy of the ledger as well as the world state (i.e., a database that contains the ledger state). Orderer nodes, on the other hand, are responsible for ordering transactions endorsed

by peers in blocks and enforcing access control for channels. In the FeatureCloud network, the "majority endorsement" policy is used. Also, raft was used as a consensus mechanism for the transaction ordering among orderers, following a "leader and follower model, in which a leader is dynamically elected among the orderer nodes in a channel. Peers and orderers have identities assigned to them by a certificate authority (CA) within the organisation they belong to (e.g., X.509). For a better distribution, each participant will normally have its own certificate authority, and the CAs of the different participants combined constitute the chain of trust.

As aforementioned, in FeatureCloud, two particular smart contracts are implemented; i) a contract for committing ML study results to blockchain in the form of a hash. The hash aggregates the hashes of all input data, output model, the ML study and the used consents for an iteration. ii) the second smart contract is used for managing commitments to operations on patient consents. This means, whenever there is a creation, update or revocation of a patient's consent to use his data, the corresponding commitment is added to the blockchain. Only authorised entities have access to the smart contract functionalities, and a participant cannot invoke the smart contract to apply state changes to assets of another participant. This means that a hospital A for example cannot apply consent updates on consents of another hospital B (i.e., the commitments on consents and not the plain text consent).

As of now, a patient may be given access to solely read and monitor the commitments of his consents on-chain. In an advanced design (see deliverable D6.4), patients may be given write access to their records on-chain (i.e., commitments). However, a pragmatic solution requires the corresponding hospital to manage the consent commitments on-chain on behalf of the patients, but the latter would still have the role of monitors, where they can check that their off-chain operations were indeed committed.

6.2 FeatureCloud Blockchain and GDPR compliance

In the following, we will assess to which extent our solution design complies with the GDPR as to whether the design enables the right to delete sensitive information from the blockchain. The assessment will be based on the rules and criteria extracted from GDPR and introduced in previous Section 5, such as who are the data controllers, whether it is possible to re-identify subjects based on blockchain data or identifiers, and whether the solution enables data erasure.

6.2.1 Personal Data and Data processing in FeatureCloud

In the context of *FeatureCloud*, patient data and consents clearly constitute personal data that identify a patient (the data subject), and their use for clinical trials or, generally speaking, for research purposes by the participants, e.g., hospitals (i.e., in this use case represent the data controllers and processors) cannot be classified as personal or household activity. As discussed in Section 5.2, Article 4(15) of the GDPR defines **data concerning health** as personal data related to the physical or mental health of a natural person, including the provision of health care services, which reveal information about his or her health status. This personal data (the raw data) is however stored locally off-chain and never leaves the participant site. In compliance with Articles 4, 6, 7 and 9 GDPR, in FeatureCloud, an informal consent is collected from patients in written or digital form if required to use their data for all or a specific type of studies. The purpose of the processing must be clearly and unambiguously specified by FeatureCloud coordinator and checked by FeatureCloud participants. Note that the consent itself represents personal data as it enables to identify the patient that gave the consent, i.e., a consent contains information such as the name, address, and the scope of the consent. Again, in FeatureCloud, obviously the consents as well are stored off-chain (in digital or paper form).

Data on the ledger in FeatureCloud and do they qualify as personal data?

In FeatureCloud, we have two different types of data that are stored on-chain: **i)** the commitments to the ML studies, and **ii)** the commitments to the consent operations.

With respect to **i)** the hash that is stored corresponds to the Merkle tree root of all input data, output models, docker configuration parameters, and the ML study, and the Merkle tree root of all input consents as follows:

```
[MLid] [Result_id] [Merkletree root of input consent] [result hash]
Example:ml0001 r0001
0EC9113091A3F7B91DEA5DAAFBBBC12DE8BCE39D7E4D55AA71A555D62639FB799
B30B8364414C181B896DE41E1B7FBB04BA639F69B115AC5E6B0BB8C22885E02B
```

Each ML study has an identifier, defined by the coordinator. Only participants involved in a given study can commit results to the study, and depending on the number of iterations, a participant may submit multiple results commitments (therefore the use of *result_id*). Note that it is possible to also aggregate the consents hashes with the ML results' hashes. Keeping consents in a separate MerkleTree, makes the proof, by the participant, that a given consent was used for a specific study shorter. Note that all raw data may be encrypted before being hashed and aggregated, and all hashes are salted or peppered.

Given that the domain space for the input data is very large, includes data in different formats, and data is aggregated (inputs, outputs and configuration parameters), then it is very unlikely for the data to be inferred from the commitment. Indeed, brute forcing this content-based re-identification may be considered “reasonably unlikely” as the time frame, the number of attempts and cost required to re-identify the subjects or the original data from this hash would be extremely high and practically infeasible. Additionally, for this scenario, it is the participants that have the responsibility to interact with the smart contracts for the ML study results, and therefore, there is no risk of identifying patients through their public keys (patients public keys are not revealed as no interactions happen between patients and the ML smart contract).

Regarding **ii)** the data that is stored on blockchain is a commitment to either a creation, update or revocation of a patient consent as follows:

```
[patient id] [consent id] [consent data Hash]
Example: p0001 c0001
60C88968F789C29610DDC920F13D7FCAAFFDD4277AD3D14F6A08504C5FB2B491
```

The transaction data are not explicitly related to a natural person but to an identifier (patient id). Therefore, to decide whether the transaction data on the FeatureCloud blockchain are personal data, we have to assess whether these identifiers (pid, cid) or consent hashes are reasonably likely to be resolved to the corresponding patient. Note that consent revocation is executed through the operation update (updating old consent hash with the hash of the revoked consent), so that the blockchain network cannot know that a consent is being revoked.

A patient id (pid) is indeed a local identifier that is only known to the participant that registered the patient and is different to the global identifier issued by for example an insurance company or ministry of health (see deliverable 6.4 for more details). Although in the example the id looks simply, in reality it can take many other forms such as the hash of a public key or a decentralised identifier (DiD). This identifier is only used locally for that specific participant, which means that if the patient is registered within multiple hospitals, each of them will have a completely different identifier for that same patient, and the identifiers are unlinkable (i.e., each hospital has its own and separate identification system). In the case where the hospital commits changes to patients' consents on their behalf, then it is the public key of the hospital which is used to submit the transactions. However, in the case where the patients themselves are enrolled as users to issue or update themselves their consents' commitments, then the public keys of the patients are used to submit transactions. In the system, it

is also possible to create a different identifier per patient and per consent operation (single-use identifiers). As such, the network will not be able to link separate consents to the same patient.

As discussed in Sections 4.2 and 4.3, these identifiers are considered as pseudonymized personal data although the re-identification of patients from the identifiers, without access to the local database of the participant, is “reasonably unlikely”. Furthermore, the re-identification of the patient through the hashed name or insurance ID is unlikely as the local identifier is generated randomly and separate from any personal data. Additionally, as the raw data of the consent is kept offline and only a salted or peppered hash of the consent is committed, then content-based-re-identification is also reasonably unlikely as the input space for consents is large and used formats differ from one hospital to another. Hashes’ salt or pepper are to be kept in a locally secure storage.

Similar to access rules defined for the ML smart contract, the functions of the consent management contract have restricted access. In the case where the participant is responsible for committing to consent operations, the public key or certificate used to do so belongs to the participant and has no link with patients. However, in the case where each patient is enrolled to manage its own consents, then patients’ certificates that allow them to interact with the blockchain become personal data. In such a case, a certificate issued by one participant can only be used for consents to that participant (i.e., to prevent linkability). Similar to the patient id, if a different certificate is used for each transaction, then re-identification may be unlikely.

In summary, no data with meaningful content is written to the blockchain. Transaction data only contains hashed data and pseudonymous identifiers. As the off-chain data remains locally within each participant, and the input space is very large, then identifier or content-based re-identification through brute forcing is reasonably unlikely. Employing a random blinding factor that is stored securely off-chain adds another level of protection against deanonymization.

6.2.2 Data Controllers, and Processors In FeatureCloud Blockchain

As aforementioned in Section 5.5, data controllers and processors play a vital role in ensuring compliance with obligations under GDPR. They must implement the appropriate measures in line with GDPR requirements. In this section, we will briefly discuss the different roles and responsibilities from a data protection perspective with a specific focus on blockchain operations. However, a complete and more detailed description of data protection roles in FeatureCloud will be discussed in future deliverable **D8.7**.

For each data processing operation, the relevant controller(s) has to be determined. In the case of FeatureCloud, the main activity is related to federated machine learning of patient data. The ML algorithm is provided by the coordinator (researcher), and the application execution is run separately by each participant (e.g., hospital). The FeatureCloud system enables synchronisation between participants, and provides an app store i) in which developers publish and certify their apps, and ii) from which coordinators select apps and construct a workflow for specific purpose.

By constructing a workflow, selecting the necessary AI apps, specifying the type of patient data required for the study, and the lawfulness for doing so, the coordinator clearly acts as a data controller.

Each participant in FeatureCloud is a healthcare data provider responsible for collecting the data, storing them and locally running the apps as instructed by the coordinator (controller). The latter might indicate that participants act as data processors but the aforementioned responsibilities speak in favour of considering them as controllers or rather joint controllers together with the coordinators. With regard to clinical trials, legal, ethical and practical arguments in favour of considering participants as data controllers are discussed in (Van Quathem, 2019).

Regarding the main activity in FeatureCloud, the federated machine learning activities on medical data, the question whether participants are always to be considered as controllers or there might be

cases where they are in the role of a processor, is to be decided by the very details of the individual case or case group. For example, it is possible to imagine a scenario where the participants collaboratively take part in determining the purpose and means of a specific study that would benefit them, e.g., to improve their treatment processes. Additionally, it may be argued that as the patient consents, in some use cases, are collected by participants with a scope that is not specific to a single study (before the study announcement), that they may qualify as data controllers (joint controllers in this case). **Deliverable D8.7** will go into more detail on these questions while in the context of this deliverable the focus lies not on the data processing (roles) regarding the federated ML but rather on the data processing regarding the blockchain.

With respect to the use of blockchain technology as an audit trail for ML studies and managing patient consents, we argue that participants qualify as data controllers, because they determine the purpose and means of using blockchain to store and process pseudonymized personal data (in this case hashes of the raw data). The mere use of the FeatureCloud blockchain infrastructure (even if a participant does not run its own peer or orderer nodes) may be considered as an implicit determination of the means of processing.

FeatureCloud utilises a permissioned private blockchain, in which peer and orderer nodes are run by entities that form a consortium, and whose identities are known and legally binding (accredited through the trust framework). These consortium entities indeed influence the determination of the means of processing as they participate in the governance and development of the network, thereby possibly qualifying as data controllers. On the other hand, one can assume that the blockchain technology serves solely as an infrastructure that anchors distributed applications (smart contracts), which then determine the means and purpose of personal data processing. As such, the applications developers may qualify as the (joint) data controllers.

In FeatureCloud, obviously the consortium members qualify as the joint-data controllers, as they collectively use the blockchain infrastructure for their own purpose. According to the CNIL⁷, when a group of participants decide to carry out processing operations with a common purpose, it is recommended that they designate a participant or an association to act as a data controller, otherwise they will all be considered joint controllers under Article 26 of the GDPR.

The smart contracts developers may be considered as solution providers. However, if they participate in the processing, then they may be qualified as data processors or data controllers depending on whether or not they are involved in the purpose of the processing. In this case, the smart contracts for managing consents and ML studies are provided by the FeatureCloud system, which then may qualify as a data processor or joint data controller depending on its role in determining the purpose. However, any future changes to the smart contracts follow a majority vote by the FeatureCloud consortium.

6.2.3 How does FeatureCloud Blockchain meet the right to rectification and erasure

Right to rectification

In FeatureCloud, patient data and consents within one participant are stored off-chain in secure local databases at the participant site. Rectifications to inaccurate or incomplete personal data in accordance with Article 16 GDPR, requires data controllers to update the data without undue delays if no exemptions apply (e.g., request manifestly unfounded or excessive).

Rectifications to patients' consents, (e.g., misspellings in the consents, inaccuracy with respect to the scope, or providing supplementary data) can still be handled off-chain through current legacy systems. However, the commitments to the new rectified consents must be stored on-chain. This means that patients can still exercise their rights to rectify their consents through the corresponding

⁷ https://www.cnil.fr/sites/default/files/atoms/files/blockchain_en.pdf

participant i.e., data controller (e.g., hospital), by providing an updated and signed paper or digital document containing the rectifications to the initial consent. If the request is deemed eligible, the rectification operation will be committed to the blockchain, and the corresponding transaction including the updated hash value will be included in a new block.

Although the commitment to the initial consent would still appear on the blockchain, it will be considered as outdated. Indeed, when creating a new consent identifier, the consent contract is designed in a way that all subsequent update operations on that identifier are chained, and the world state (blockchain state) would contain the last value (commitment) to that consent ID. Only authorised entities have the access rights to apply updates on that specific consent. Whether it is the patient or the participant on its behalf that submits consent commitment transactions to the blockchain, a patient always has the possibility to monitor/check that the operation took place on-chain.

Regarding the debate on whether adding the rectified data on-chain in a subsequent block without deleting the old inaccurate/incomplete data satisfies Article 16 of the GDPR, FeatureCloud keeps the original data off-chain. The inaccurate data can therefore be deleted, but the corresponding hash would still remain on-chain. However, it will be reasonably unlikely to obtain that inaccurate data from the hash (cf. Section 6.2.1). Also, querying the blockchain for the patient's consent will return the last corresponding commitment from the World state (a database that contains the blockchain state).

Rectifications to other patient personal data (e.g., treatments or diagnoses) follow the same off-chain process, except that they are not committed on-chain (out of scope of FeatureCloud). In FeatureCloud, only input data used for FL studies are committed in conjunction with other data. As the ML studies input data and output models are only committed on-chain after processing, then they are not affected by the future rectification's requests (the requests come after the processing). Personal data rectification requests received before the processing can be taken into consideration, and will not have impact on the on-chain data, as the commitment only happens after processing.

Right to erasure

Similar to the right to rectification, GDPR provides data subjects the right to request the deletion of their personal data, and the controller shall have the obligation to erase and stop processing them without undue delay if no exemptions apply. While Article 17(1) GDPR lists a number of grounds to exercise the right to erasure, Article 17(3) provides exemptions where the obligation to apply this right by the controllers no longer holds. Particularly relevant here is the right to erasure can be exercised if “the personal data are no longer necessary in relation to the purposes for which they were collected or otherwise processed” (Article 17(1)(a) GDPR) or if “the data subject withdraws consent on which the processing is based” (Article 17(1)(b) GDPR). One particular exempt from the right to erasure that may eventually be used in the context of FeatureCloud is: “for archiving purposes in the public interest, scientific or historical research purposes or statistical purposes in accordance with Article 89(1) in so far as the right referred to in paragraph 1 is likely to render impossible or seriously impair the achievement of the objectives of that processing;”

This is due to the fact that the processing purpose in FeatureCloud, with respect to blockchain, is the documentation or archiving of consents (commitments) and ML studies contexts (also commitments) to ensure the integrity and auditability of used data and corresponding consents. Deleting the personal data would render the processing (auditability and integrity checks) impossible (as it will not be possible to replay the ML algorithms, but also to check whether a consent was given at a certain point in time).

In the case where no exemptions apply (e.g., the FeatureCloud audit already took place, and data is no longer needed for audit purposes), possible solutions to erase personal data under the obligation of GDPR include the complete destruction of the data, or its anonymization following the

discussions in Section 5.2 and 6.2.1. Note that truly anonymized data is not subject to GDPR. As discussed in previous sections, the right to erasure and data deletion may at first look contradictory to one of the main properties of blockchain, i.e., that is immutability (data stored on blockchain is meant to leave forever).

Complying with Article 17 GDPR if no exemption applies requires all personal data to be removed from all peers (nodes), which hold a copy of the ledger. This means that upon the reception of an eligible erasure request by a participant, the latter not only has to erase the data from its own local databases/servers, but also needs to initiate erasure from all joint-controllers (i.e., the blockchain nodes). In a permissionless setting, this might be technically infeasible, given the dynamicity of the network nodes and the uncertainty about who is accountable in such a setting. Furthermore, in some blockchain designs such as Bitcoin, nodes belong to different jurisdictions in which GDPR may not have extraterritorial overreach. In FeatureCloud, however, the blockchain is permissioned where only entities accredited by the trust framework are allowed to operate nodes. Identities of all nodes are known to the consortium and legally binding. This means that the data controllers and processors are known in advance and adhere to compliance with GDPR requirements during onboarding.

One suggested solution for erasing data on blockchain is to use pruning, where a part of the blockchain is pruned after a specific retention period. This means that all blocks exceeding that time frame will be deleted. We think this is not a viable solution in our setting as it will no longer be possible to check whether a consent was given at a certain point in time where the study took place. This is because ML studies use data inputs whose consents are given at different points in time. To explain this, let's assume that the time frame to retain the data is 5 years, and that two different patients' data as well as their corresponding consents are 4 and 2 years old respectively. Now assume that both data are being used for a new study, which is then committed. As all blocks older than 5 years get removed, then in the following year after the study, the block containing the first consent commitment will also be 5 years old and thus removed. The dependence of the study, which is then only 1 year old, on that first consent, will later make it non-auditable. Moreover, pruning will render the verification of the world state in certain cases impossible as it requires the full history of state changes.

Another solution mentioned in the literature consists in forking the blockchain, where the new fork will not include the block containing the data to be removed. The fork would require a majority vote by the nodes participating in the network. In addition to similar arguments against pruning, we think that this solution is impractical and a continuous number of erasure requests would exponentially increase the number of branches as well as costs and efforts to maintain a valid state among nodes. This will also create indecision and uncertainty about which branch to follow.

To comply with the right to erasure, FeatureCloud instead follows a methodology based on data protection by design and default as prescribed by Article 25 GDPR⁸ (“implement appropriate technical and organisational measures, such as pseudonymisation, which are designed to implement data-protection principles, such as data minimisation, in an effective manner and to integrate the necessary safeguards into the processing in order to meet the requirements of this Regulation and protect the rights of data subjects”). The legal requirements collected in deliverable D6.3, allowed us to consider data protection rights while designing the FeatureCloud solution for managing consents and ML studies. By using private permissioned blockchain, FeatureCloud enforces data and smart contract access controls as well as stricter usage of the rules. In particular, it limits access to who can read specific data or who can act on existing data (such updating consents).

In FeatureCloud, blockchain technology is mainly used to certify on-chain data that is stored off-chain. This means that patient data records as well as the actual consents are never stored on the blockchain. What is stored on-chain is a local identifier (only known to the respective hospital), and

⁸ <https://gdpr-info.eu/art-25-gdpr/>

a hash of the content (cf. Section 6.2.1). The mapping between the patient local identifier (on-chain) and the patient identity (global identifier) is stored locally off-chain in a secure location. Similarly, the original consent is also stored off-chain with the blinding factor used for the hashing function. Therefore, erasing the off-chain data will render the respective on-chain data obsolete and hence non-personal. By deleting the off-chain data (the data and the mappings), it will be reasonably unlikely to re-identify the patients based on the identifiers or the content hashes stored on-chain. This is particularly due to the domain space being sufficiently large in addition to the blinding factor used for hashing, thus rendering brute forcing significantly hard. An additional protection measure is to also encrypt the data (consents) with a secure encryption scheme before hashing. By deliberately deleting or overwriting the encryption keys as well as the off-chain data (i.e., crypto-shredding), it will no longer be possible to reconstruct the data from the hashed encrypted data. In addition, it is also possible to use a separate identifier (for the same patient) each time a new consent is created (cf. Section 6.2.1).

In most cases, it is the participant (on behalf of the patient) who interacts with blockchain for committing consents (if patients cannot use cryptographic tools). However, in the case where patients themselves interact with the blockchain for reading their data commitments (or if they are given write access), then it is recommended that they obfuscate their IP addresses through the use of anonymisation tools such as Tor or VPN. Although IP addresses are not logged by the peers in FeatureCloud, utilising anonymization tools prevents correlating users IP addresses with their consent commitments stored on the blockchain. Furthermore, as the FeatureCloud blockchain setting is permissioned, then all peers within the consortium undergo an onboarding process in which legal responsibilities are assigned to them with which they must comply. This includes the deletion of personal data collected during the interaction with the blockchain if requested by the users (eventually the IP addresses).

Additionally, in FeatureCloud, participants have the possibility to randomly and periodically push consent updates to the blockchain to hide meta-information about patients' consents such as how many and when consents were given by a patient. This means that a participant can periodically invoke Create or Update consent operations with hashes of random nonces (instead of consent hashes). This makes it harder for an outsider (other peers) to infer some metrics such as how many consents a hospital has, or how many consents are given by a patient ID. For the other peers, it will not be distinguishable whether or not the transaction corresponds to a real consent. For the auditors, later it can be easily proven which commitments correspond to actual consents.

Similar to consent commitments, ML studies related transactions do not include any identifiable personal data as discussed in Section 6.2.1. This is because the hash that is stored on blockchain corresponds to the Merkle Tree root of all input data, output models, the ML study and the configuration parameters, thereby making the domain space so large that is unlikely to be brute forced. Note that the Merkle Tree root of all consents used for the study is committed as a separate parameter. The reason is to make the proof that a patient consent was used for a specific study shorter and faster (proofs are provided by participants to either patients or auditors). Note that the blinding factors used for consents within the consent management smart contract should be different from the ones used for consents with the ML study contract (even if it is the same consent). This is to prevent other peers from correlating consents to ML studies using analysis techniques (for example by generating Merkle Trees from existing committed consents and comparing the root to the one committed within a study e.g., brute forcing).

Private Data Collections and the right to erasure:

In Hyperledger, it is possible for a subset of peers to privately share data while hiding it from the remaining peers of the same channel, i.e., using private data collections. This has the advantage over creating separate channels that it avoids additional administrative overhead such as channel

configuration, defining policies or maintaining chain code versions. The private data is sent peer-to-peer to only authorised organisations. While this data will be privately stored at each authorised peer, a hash of the data will be stored on the ledger of each peer in the channel. When invoking the chain code, the private data is sent in a transient field to prevent it from getting stored within the transaction context as input parameter.

One of the initial objectives of private data collections is to provide a GDPR friendly solution, in which data can be kept privately in private databases, while their hashes are committed to the public chain within the channel. As the private data resides in a transient data store at each authorised peer, then it becomes possible to delete it upon request by the user (following its right to erasure). The erasure of the private data on all peers will make the corresponding hashes stored on the channel public ledger obsolete. Note that it is possible to invoke deletion before the blockToLive expires (blockToLive is an argument which specifies a block after which private data will be purged).

Although the concept is interesting, at this stage of FeatureCloud, no private data is shared among participants, i.e., data is kept locally at each participant's site. However, in the future we will investigate whether this concept can be used to facilitate data sharing and enable global discovery (in deliverable D6.6). Additionally, so far, we use a similar approach where only commitments are stored on the public ledger of the FeatureCloud channel, however the data itself is not shared with any peer. Moreover, there exist a number of issues related to the use of private data collections in Hyperledger that remain unresolved with respect to GDPR^{9 10}. For example, when invoking the deletion of private data, some of the information is kept as part of the peer logs (private history database), which might include personal information. A request for comments (RFC) is being discussed to suggest a new operation (private_data_purge), which also purges all historical data¹¹.

To summarise, the FeatureCloud data protection by design approach for complying with the right to erasure in the light of the use of Blockchain technology is that if, based on the considerations above, Art 17 GDPR requires erasure, off-chain data is deleted, which turns the corresponding on-chain data into non-personal data because without the deleted off-chain data the identification of the data subjects relating to the on-chain data becomes practically impossible.

7 Conclusion

While blockchain technology may readily provide transparency and immutability for the data recorded on a shared ledger, these very characteristics can be problematic in regard to privacy and data protection requirements such as GDPR. Indeed, the ability to rectify or delete data constitutes a necessity for compliant systems. At first, it would appear that in such cases incorporating BT is not an ideal approach. However, with careful design considerations, the advantages of BT can be leveraged while avoiding these privacy issues. In FeatureCloud, an approach that incorporates data protection by design and default was adopted. The use of permissioned blockchain for FeatureCloud enables more control over the blockchain network, and easier identification of data controllers as the entities to be held accountable for compliance with the GDPR. Furthermore, keeping the actual data off-chain and their commitments on-chain facilitates auditability and adds trust to consent management while still enabling compliance with GDPR and the right to erasure or rectification.

⁹ <https://jira.hyperledger.org/browse/FAB-5097>

¹⁰ https://docs.google.com/document/d/1O5UCzf8H_kN4iqMgHm8haDhqacsCBEWNY-dhDG7qgYc/edit#heading=h.9puzkv11asp9

¹¹ https://hyperledger.github.io/fabric-rfcs/text/0000-private_data_purge.html

8 References

- Finck, M., 2019a. Blockchain and the general data protection regulation: can distributed ledgers be squared with European data protection law? European Parliament. Directorate General for Parliamentary Research Services., LU.
- Finck, M., 2019b. Smart contracts as a form of solely automated processing under the GDPR. *Int. Data Priv. Law* 9, 78–94. <https://doi.org/10.1093/idpl/ipz004>
- Finck, M., Pallas, F., 2020. They who must not be identified—distinguishing personal from non-personal data under the GDPR. *Int. Data Priv. Law* 10, 11–36. <https://doi.org/10.1093/idpl/ipz026>
- Gauravaram, P., 2012. Security Analysis of salt||password Hashes, in: 2012 International Conference on Advanced Computer Science Applications and Technologies (ACSAT). Presented at the 2012 International Conference on Advanced Computer Science Applications and Technologies (ACSAT), pp. 25–30. <https://doi.org/10.1109/ACSAT.2012.49>
- Ghesmati, S., Fdhila, W., Weippl, E., 2022a. SoK: How private is Bitcoin? Classification and Evaluation of Bitcoin Privacy Techniques, in: Proceedings of the 17th International Conference on Availability, Reliability and Security, ARES '22. Association for Computing Machinery, New York, NY, USA, pp. 1–14. <https://doi.org/10.1145/3538969.3538971>
- Ghesmati, S., Fdhila, W., Weippl, E., 2022b. User-Perceived Privacy in Blockchain.
- Van Quathem, K., 2019. The GDPR and Clinical Trials – Are Study Sites Controllers or Processors?, in: *Gesetz Und Rech.* Presented at the pharmind, ECV • Editio Cantor Verlag, Aulendorf (Germany).

9 Table of acronyms and definitions

BT	Blockchain Technology
concentris	concentris research management GmbH
CA	Certificate Authority
DES	Data Encryption Standard
DLT	Distributed Ledger Technology
FL	Federated Learning
GDPR	General Data Protection Regulation
GND	Gnome Design SRL
IP Address	Internet Protocol Address
ISP	Internet Service Provider
ML	Machine Learning
MS	Milestone
MUG	Medizinische Universitaet Graz
Patients	In this deliverable, we use the term “patients” for all research subjects. In FeatureCloud, we will focus on patients, as this is already the most vulnerable case scenario and this is where most primary data is available to us. Admittedly, some research subjects participate in clinical trials but not as patients but as healthy individuals, usually on a voluntary basis and are therefore not dependent on the physicians who care for them. Thus, to increase readability, we simply refer to them as “patients”.
RI	Research Institute AG & Co. KG
SBA	SBA Research Gemeinnützige GmbH
SDU	Syddansk Universitet
TUM	Technische Universitaet Muenchen

UHAM	University of Hamburg
UM	Universiteit Maastricht
UMR	Philipps Universitaet Marburg
WP	Work package