






## Article

## DEFS—Data Exchange with Free Sample Protocol

Rafael Genés-Durán <sup>1,\*</sup>, Juan Hernández-Serrano <sup>1</sup>, Oscar Esparza <sup>1</sup>, Marta Bellés-Muñoz <sup>2</sup>  
and José Luis Muñoz-Tapia <sup>1</sup>

<sup>1</sup> Department of Network Engineering, Campus Nord, Universitat Politècnica de Catalunya, 08034 Barcelona, Spain; j.hernandez@upc.edu (J.H.-S.); oscar.esparza@upc.edu (O.E.); jose.luis.munoz@upc.edu (J.L.M.-T.)

<sup>2</sup> Department of Information and Communications Technology, Tànger Building, Universitat Pompeu Fabra, 08018 Barcelona, Spain; marta.belles@upf.edu

\* Correspondence: rafael.genes@upc.edu

**Abstract:** Distrust between data providers and data consumers is one of the main obstacles hampering the take-off of digital-data commerce. Data providers want to get paid for what they offer, while data consumers want to know exactly what they are paying for before actually paying for it. In this article, we present a protocol that overcomes this obstacle by building trust based on two main ideas. First, a probabilistic verification protocol, where some random samples of the real dataset are shown to buyers in order to allow them to make an assessment before committing any payment; and second, a guaranteed, protected payment process enforced with smart contracts on a public blockchain that guarantees the payment of data if and only if the provided data meet the agreed terms, and that honest players are otherwise refunded.

**Keywords:** blockchain; marketplace; privacy; fairness; conflict resolution; non-repudiation; probabilistic verification



**Citation:** Genés-Durán, R.; Hernández-Serrano, J.; Esparza, O.; Bellés-Muñoz, M.; Muñoz-Tapia, J.L. DEFS—Data Exchange with Free Sample Protocol. *Electronics* **2021**, *10*, 1455. <https://doi.org/10.3390/electronics10121455>

Academic Editor: Priyadarsi Nanda

Received: 16 May 2021

Accepted: 15 June 2021

Published: 17 June 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

The use of data has increasingly become a crucial factor in the success of businesses. Research has shown that proper use of big data techniques helps to identify new insights, optimise operating processes, and make better and faster decisions [1]. In this context, ecosystems have grown to fulfill the data needs of diverse actors, such as data suppliers, data custodians, or data aggregators. As a result, businesses not only collect and analyse the data they generate but increasingly rely on third party data to enhance their business value. The necessity of exchanging data between different parties gives rise to an ecosystem that has an inherent regulatory complexity and a need for privacy. In general, making proper data agreements is not easy, especially the task of valuing data and convincing customers of their value without giving them away [2]. The creation of marketplaces addresses many of these problems. Allowing providers and consumers to deal with common interests in a platform where both parties can meet each other and trade information solves the integration problem of connecting consumers and providers.

In this article, we focus on the problem of convincing consumers of data value, which can be seen as a form of lack of trust towards data providers. Traditionally, this problem could not be solved without previously establishing a certain confidence between parties. This fact represents an entry barrier to new providers in the market, hurting competence and thus reducing utility for consumers. Achieving the exchange of virtual products between many parties while minimising risks is the main goal of virtual commerce. In order to exchange value safely, it is essential to ensure that consumers get the product they pay for and that providers get paid. These two things are often carried out without any strict protocols and guaranteed just by existing trust. Typically, counterparties that know each other from previous experience or whose future interests are aligned are confident that no intent to scam will be made by the other party, since confidence is often more beneficial than gains from fraud.

Nonetheless, when stronger assurance than that is needed, it is a common practice to use a trusted third party (TTP) that all parties trust to guarantee that the process is carried out correctly by all involved individuals. TTPs solve the crucial aspect of minimising risks, but ensuring its viability entails an extra cost for all parties. Moreover, centralising interactions between businesses via the TTP generates a single point of failure that could produce critical delays and denial of services. Distributed ledger technologies (DLTs) can be seen as a paradigm shift when it comes to the need of TTPs. Using DLTs, all participants in the network can maintain a set of synchronised data (who owns what) without the need for a central authority (TTP) guaranteeing integrity, fairness, and data availability. In addition, recent studies have shown that replacing TTPs with DLTs represents an important optimisation of time and overall costs [3].

In this article we present DEFS, a protocol that addresses the lack-of-trust problem between providers and consumers in a data trade. Our protocol preserves the security, privacy, and fairness standards that marketplaces should guarantee, and it also includes the capability of checking some sample portions of the dataset before committing to purchase to enhance the trust of the consumers in the data value.

The article is organised the following way: In Section 2, we give an overview of the technologies used in our protocol. Section 3 contains the state of the art on decentralised data marketplaces. In Section 4, we explain our protocol. First, we give a general overview, and afterwards, we provide a detailed description of each step of the protocol. In the following section, Section 5, we present a security analysis and, finally, we conclude in Section 6.

## 2. Background

### 2.1. Distributed Ledger Technologies

The main technology to build a public ledger is a blockchain network. In a blockchain network, users can run a blockchain node to send their transactions or use some available node that allows them to do so. Then, in a distributed way, the blockchain network can create a unique sequence of ordered transactions. In more detail, the network creates a chain of blocks following a consensus algorithm to order transactions [4]. A block contains several transactions, and an important property of these systems is that once the consensus algorithm definitively accepts a block, all nodes get to know this block and it becomes impossible to manipulate or delete it [5].

In a blockchain network, users can own one or more accounts. Accounts are identified via a public identifier (usually derived from a random public key using a hash function). New blockchain accounts can be created by simply generating a pair of asymmetric keys and deriving the account identifier from the public key. In general, account identifiers are not directly linked with any user data, so they can be considered pseudo-anonymous identifiers. Transactions carry the source account identifier and a destination account identifier, and they are all digitally signed using the private key of the source account. All the nodes that form the blockchain network see the same state (also known as world state) that results from executing all the transactions in order.

### 2.2. Smart Contracts

Some blockchains not only allow executing regular transactions that modify the cryptocurrency balances on the ledger but also have the capability of deploying and executing public and auditable programs called smart contracts. Smart contracts have their own state, and in their code we can define the business logic we want to process transactions. Once a smart contract is deployed in the blockchain network, its code is replicated on every node and, consequently, these programs have the same availability and integrity as regular transactions. The Ethereum [6] mainnet is a good candidate to implement our proposal because it is a public blockchain, capable of running smart contracts, and it is the platform of choice for many developers for implementing decentralised applications (DApps).

### 2.3. Merkle Hash Trees

A Merkle hash tree (MHT) is an authenticated data structure where every leaf node of the tree contains the cryptographic hash of a data block and every non-leaf node contains the concatenated hashes of its child nodes [7]. MHTs allow linking a set of data to a unique hash value, the Merkle hash tree root (MR), allowing efficient and secure verification of the consistency and content of large sets of data.

Figure 1 contains an example of a MHT with 8 leaves. To show that a certain value is stored in a leaf of the MHT, one can create a Merkle proof (MP), which consists of a list of the additional nodes required to compute the root of the tree. For instance, a MP showing that  $h_3$  is stored in the MHT from Figure 2 would consist of the nodes

$$MP(h_3) = \{h_2, h_{01}, h_{4567}, h_{01234567}\}.$$

Note that with  $h_3$  and the first three nodes of this list, anyone can compute the root of the tree. If the root matches  $h_{01234567}$ , then the proof is valid proof of membership for  $h_3$  in the tree.

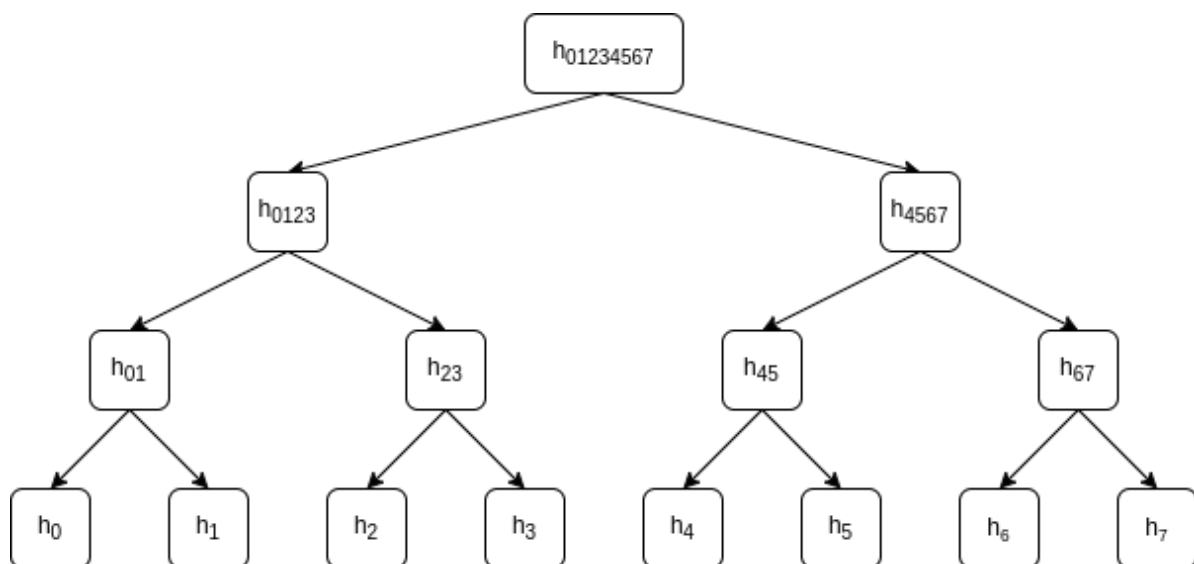


Figure 1. Merkle hash tree (MHT) of 8 leaves.

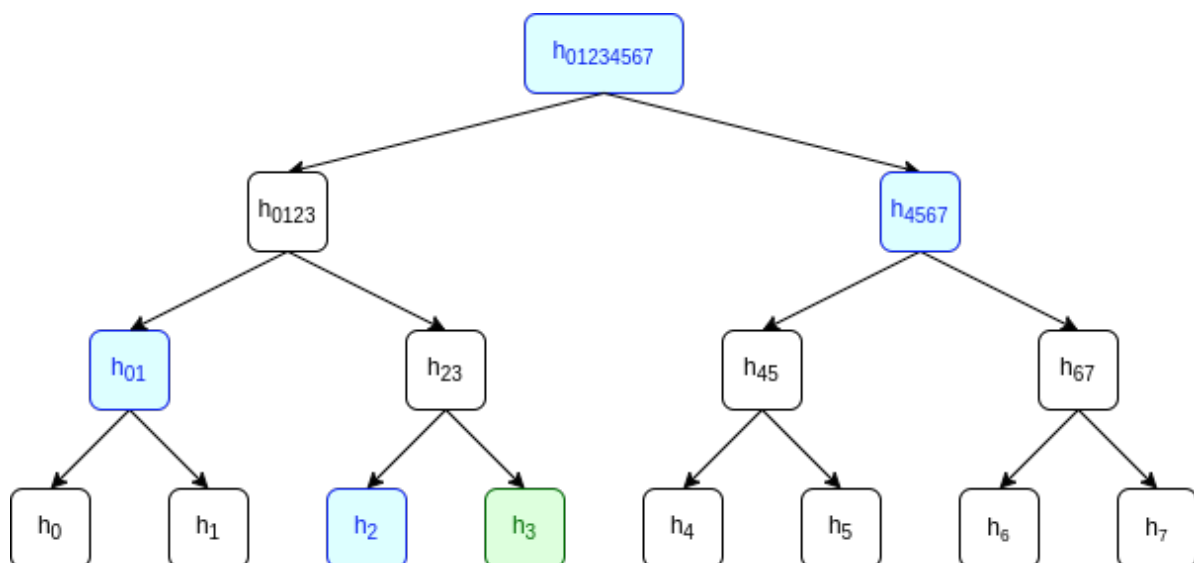


Figure 2. The Merkle proof (MP) of  $h_3$  is the set  $\{h_{01234567}, h_{4567}, h_{01}, h_2\}$ , which contains the nodes needed to compute the MR of the tree.

The security of a MP reduces to the collision resistance of the underlying hash function [8]. For this reason, we assume the hash function  $H$  used to build MHTs is cryptographically secure. That is, that the probability of finding a preimage or a hash collision is negligible [9].

### 3. State of the Art

Traditional data marketplaces build trust by requiring identifying information to the different stakeholders. This information and the process to get it is typically known as know your customer (KYC). The belief is that KYC constitutes an entry barrier to fraudsters, since KYC data would help victims to refer to law enforcement if tangible fraud is committed. However, marketplaces with KYC lead to regulatory complexity and are, in general, difficult to operate with. This fact hinders the growth of digital-data trading and can make stakeholders feel that their privacy has been violated. In addition, law execution is slow and, in some cases, might be even useless. For the previous reasons, there have been attempts in the research community to design new approaches to mitigate these limitations with technology and not regulations. In this context, decentralised marketplaces have arisen as a solution to enhance security, sovereignty and trust in data exchanges [10–12].

One interesting initiative is GAIA-X [13], which is an European project created to develop the foundations for a federated open-data infrastructure connecting both classical architectures with decentralised infrastructures in order to build a transparent ecosystem for the end users taking advantage of the decentralised benefits.

One of the main technologies that is fostering data marketplaces is the Internet of Things (IoT), which generates huge amounts of data from sensors and devices. The increasing necessity of monetising these data is also pushing research. In the literature, we can find several works that propose decentralised marketplaces for IoT using distributed ledger technologies to enhance the data exchanges with transparency, trust, and integrity [14–16]. Among others, decentralised marketplaces are being implemented in new disruptive scenarios such as artificial intelligence [17], smart cities [18,19], and connected cars [20]. In fact, the value of data is becoming more important in business interactions, which is reflected in the new technologies and their necessity to generate this new era of decentralised marketplaces.

A remarkable example of a decentralised data trading solution is presented in [21]. As in our protocol, the data on sale are not stored on the blockchain but in some external (and possibly distributed) storage platform. Similar to our protocol, the proposed solution symmetrically encrypts data on sale and uses a MHT of cryptograms to register the associated trades on the blockchain. However, the solution proposed not only requires to generate symmetric cryptograms but also the need of asymmetrically sign each of these cryptograms. Additionally, authors propose using plaintext checkable encryption (PCE) [22] to check, on chain, that cryptograms have been correctly encrypted. In our protocol, we avoid using asymmetric encryption, which is much slower than symmetric encryption. With DEFS, we achieve a faster and easier solution by providing structure to symmetric keys and generating a MHT with these keys to allow solving disputes with regards to data encryption. In addition, Ref. [21] considers three roles: data buyers, data sellers, and miners. The main problem of directly involving miners in the implementation of the solution is that the mining software then needs to be modified, which is, in general, not a trivial thing to do. Mining software is extremely subtle, since any error in an implementation can lead to a lack of consensus in the network. In our protocol, we also consider the roles of data buyers and data sellers, but the role of miners is abstracted, and we use the API provided by smart contracts which is a much easier and safer way of implementing the logic of data trades in the blockchain.

Another remarkable implementation of a decentralised data trading solution is presented in [23], where the authors present SDTE, a secure blockchain-based data trading ecosystem. As with our protocol, SDTE tries to mitigate the existence of dishonest parties in data exchanges. However, SDTE focuses on an scenario in which the buyer does not need to have access to a complete dataset but it only needs the findings from the data analysis. For

this case, SDTE proposes a data processing-as-a-service, where the buyer is paying for the analysis of the seller's dataset. SDTE is build using an Intel's SGX-based secure execution environment to protect the data processing, the source data, and the analysis results. As we will show in the following section, DEFS is not designed as data processing-as-a-service but as data exchange-as-a-service. In the latter, the seller wants to buy the complete dataset, not computed data. For this scenario, DEFS provides a probabilistic verification protocol and a conflict resolution protocol that is guaranteed and supported by a smart contract.

#### 4. Data Exchange Protocol

In this section we introduce DEFS, a protocol that addresses the problem of data trading between a provider and potential consumers using a smart contract deployed in the blockchain as a broker. To mitigate gender issues when referring to a single provider and a single consumer, we will assume the provider is a woman and the consumer a man.

As we explained before, the use of DLTs can replace the role of TTPs in payment processes. When using DLTs, participants in the network can maintain synchronised data and share payment information without the need for a central authority, in this way guaranteeing the integrity, fairness, and availability of the data. In this manner, DEFS makes use of a smart contract to preserve the security and privacy standards that marketplaces should guarantee.

Another gap to cover in this data trading scenario is generating trust between data consumers and data providers. Here comes the novelty of DEFS: our proposed data exchange protocol is designed with the capability of checking random samples from the dataset so that consumers are able to infer if the complete dataset is worth paying for, enhancing the trust from the consumer's side. On the other side, the smart contract acts as a broker during the payment procedure, ensuring providers that they will receive the payment for the data they exchanged.

##### 4.1. Protocol Overview

First, we give a general overview of DEFS, and in Section 4.3, we describe in greater detail all the steps the entities involved (consumer, provider, and smart contract) should follow. We assume that before starting the protocol, a data provider advertises their data to the public using off-blockchain means, such as a data marketplace. Then, a consumer interested in a particular dataset contacts the provider, who starts the DEFS protocol to perform the data exchange and payment. To prevent potential extensive leaks of the data, it is important that the DEFS protocol is executed independently per each individual consumer. DEFS consists of three different phases:

1. **Protocol preparation.** In this initial phase, the provider prepares not only the data to be exchanged but also all the parameters and cryptographic material necessary to demonstrate that the data exchange is secure and private. More specifically, the provider:
  - Divides the complete dataset in portions. These portions are chosen randomly (not consecutively) from the dataset.
  - Generates a seed to generate symmetric cryptographic keys.
  - Uses these keys to create a MHT, whose root can be used to check the correctness of this cryptographic material.
  - Encrypts a random permutation of the data portions with the keys, obtaining an encrypted and randomised version of the whole dataset.
  - Creates another MHT using the hashes of these cryptograms as leaves, whose root can be used to verify the correctness of the cryptograms generated.
  - Deploys a smart contract in the blockchain that includes among other information, the roots of the previous trees.

If a consumer has interest in obtaining the dataset, the protocol continues as follows:

- The consumer receives the whole dataset encrypted but it cannot be decrypted at that very moment.
- The consumer queries the smart contract to obtain the root of the tree of cryptograms and verifies that all the cryptograms belong to this tree.

As previously stated, this is only a brief summary of the steps to follow in this phase. A more exhaustive explanation of the protocol preparation phase can be found in Section 4.3.3. At this point, all entities (consumer, provider, and smart contract) are ready to start the protocol execution phase, in which the consumer will have access to the complete dataset and will perform the payment.

2. **Protocol execution.** In this phase, the consumer gets some samples of the dataset (for free) to evaluate if it is worth to pay for the whole set, and if so, he will obtain the dataset and the provider will get paid:
  - The consumer chooses at random some sample portions to be revealed. Note that the provider committed the shuffled encrypted data at the very beginning of the protocol. Since the consumer requests random samples, neither consumers nor providers have control over the samples that will be revealed.
  - The provider discloses the keys for those samples so the consumer can evaluate the quality of the dataset.
  - If the consumer is not convinced, the protocol ends here. However, if they decide that it is worth paying for the dataset, they commit the payment to the smart contract.
  - The provider is asked to publish the seed (that will disclose all the encryption keys) in the smart contract.
  - If the consumer is able to properly decrypt the dataset, after a timeout, the provider gets paid and the protocol ends.
  - If the consumer is able to prove that there were problems with the previous procedure, he starts a conflict resolution phase to obtain a refund.

A more exhaustive explanation of the protocol execution phase can be found in Section 4.3.4.

The following phase will only be needed in case the consumer considers that he has been cheated on.

3. **Conflict resolution.** This phase is optional and only takes place if the consumer detects provider misbehaviour. The conflict resolution can end with a refund if the consumer is able to demonstrate one of the following misbehaviours:
  - A key is not properly generated.
  - A cryptogram does not have the proper format when decrypted.

A more exhaustive explanation can be found in Section 4.3.5.

#### 4.2. Protocol Properties

The main properties provided by our protocol are as follows:

1. **Data sample evaluation.** The consumer gets a free set of fair samples of the data being traded before paying. The protocol ensures that neither the consumer nor the provider are able to manipulate the chosen data or select specific samples.
2. **Payment guarantees.** The provider gets paid if and only if the consumer has access to the whole set of data. That is, the consumer cannot get the data without paying for it and the provider does not get paid without disclosing the data.
3. **The solution is cost-efficient.** Due to high fees on public ledgers, DEFS minimises the amount of data stored on the ledger, which is also independent of the quantity of data traded. This way, both the amount of data stored and the number of interactions with the distributed ledger are constant.



4. **Non-repudiation.** The DEFS protocol ensures that any party involved in the exchange is not able to cancel and/or deny the data exchange once an agreement is made.
5. **Liveness.** The different timeouts guarantee that the protocol reaches a final state, even when one of the parties quits in advance.

#### 4.3. The DEFS Protocol

In this section, we describe the DEFS protocol. We establish the notation in Section 4.3.2. The procedure before initiating the exchange is detailed in Section 4.3.3. Then, the interactions to do a fair transaction are explained in Section 4.3.4. The conflict resolution is detailed in Section 4.3.5. Finally, we present the state diagram of the smart contract in Section 4.3.6.

##### 4.3.1. Requirements

The DEFS protocol assumes that measures to meet the following requirements are already in place:

**Secure off-chain channel between provider and consumer:** It is assumed that the off-chain channel between consumer and provider is end-to-end protected. This requirement can be easily met by using the widely supported TLS protocol—e.g., with HTTPS. TLS only requires the server to hold a valid certificate (and its complementary private key) in order to create the secure channel. That is to say, consumers just need a valid TLS client, which is implemented by default in most programming languages, application frameworks, and/or web browsers.

**Validation of data blocks' format:** It is assumed that consumers can verify received data blocks according to a previously agreed schema. The process that verifies that a data portion meets a predefined format is usually called a validator. Validators are used by many technologies to check received responses before processing them. In object-oriented programming, this process is usually done by trying to parse the response as a given type of object, which will produce an error if it does not. There are also specific standards with well-known implementations, such as JSON-LD [24], that help define and validate specific data schemas.

**Identification system:** The DEFS protocol releases random samples of the dataset to potential consumers before they commit paying for the entire dataset. However, there is a risk of an attacker using multiple identities to retrieve a representative portion of the dataset for free. In this context, DEFS assumes that there are off-chain solutions run by the providers which could effectively limit the amount of identities an attacker could take. A known example is binding the identity to an e-mail account or a mobile phone. The provider should decide the most suitable authentication method depending on the price of the traded data and the type of consumers. For example, in some cases, authenticating with an e-mail can be enough. In other scenarios, e-mail might not be enough because it is not hard to generate multiple “identities” based on different e-mail accounts. In the latter case, providers might require authenticating with a mobile phone or even with both factors. In some specific cases, authentication could involve more factors, such as physical key generators, smart cards, etc.

##### 4.3.2. Notation

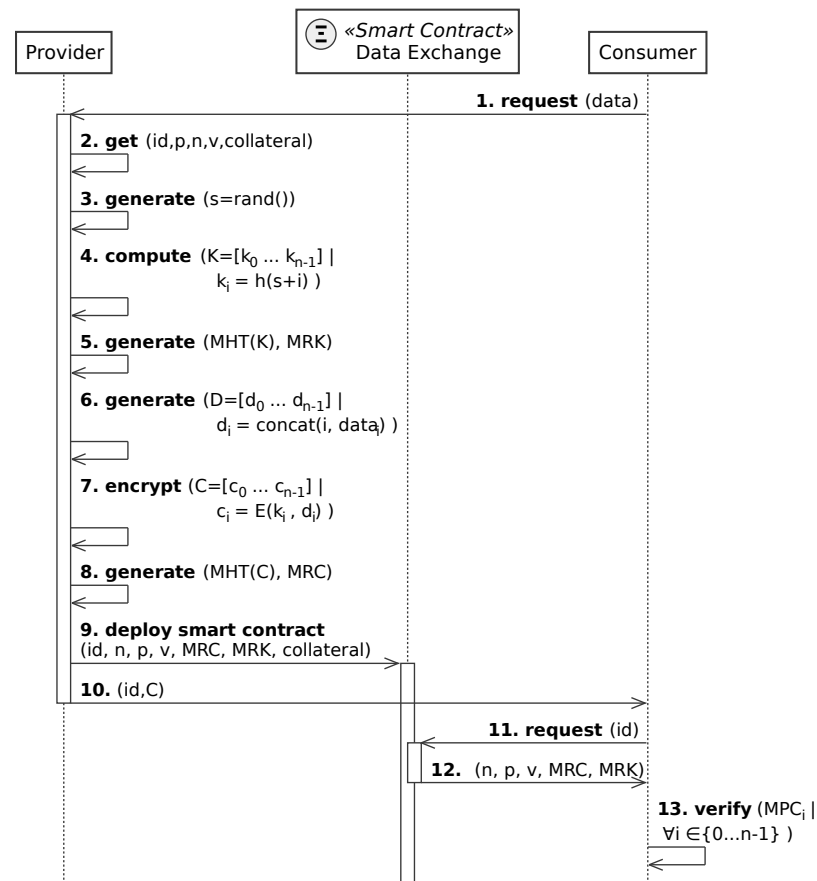
The notation of the DEFS protocol is summarised in Table 1.

##### 4.3.3. Protocol Preparation

In Figure 3, we detail the interactions between provider, consumer, and smart contract during the preparation phase of the DEFS protocol.

**Table 1.** Notation for the DEFS protocol.

Notation	Description
<i>collateral</i>	Price that the provider must pay in order to ensure fairness.
$C = \{c_0 \dots c_{n-1}\}$	Encrypted data portions (cryptograms) such that $c_i = E_{K_i}(d_i)$ .
$D = \{d_0 \dots d_{n-1}\}$	Data portions.
<i>e</i>	Index of an invalid cryptogram.
<i>id</i>	Data-exchange identifier.
$K = \{k_0 \dots k_{n-1}\}$	Encryption keys for data portions, such that $k_i = \text{hash}(s + i)$ .
$MHT(C)$	Merkle hash tree of cryptograms.
$MHT(K)$	Merkle hash tree of keys.
<i>MRC</i>	Root of the Merkle hash tree of cryptograms.
<i>MRK</i>	Root of the Merkle hash tree of keys.
$MPC_i$	Merkle proof of a cryptogram with index <i>i</i> .
$MPK_i$	Merkle proof of an encryption key with index <i>i</i> .
<i>n</i>	Number of data portions.
<i>p</i>	Price of the dataset.
$R = \{r_0 \dots r_{v-1}\}$	Set of indexes of the sample portions to be revealed.
<i>s</i>	Seed. Random number for key generation.
<i>v</i>	Number of sample portions to be revealed.

**Figure 3.** Protocol preparation: sequence diagram.



The steps of this phase are enumerated and explained immediately below:

1. **Consumer**→**Provider**: Request *data*.  
The protocol starts with the consumer's interest in a dataset. Through the marketplace, the consumer requests some offered *data* to the provider. Note that each time the consumer desires a dataset, a new instance of the DEFS protocol is required.
2. **Provider**: Set *id, p, n, v, collateral*.  
In this step the provider has to decide the main parameters associated to the dataset. These parameters are the identifier of the data exchange (*id*); the price of the dataset (*p*); the number of portions in which the dataset will be divided (*n*); the number of sample portions to be revealed before the payment (*v*); and, finally, the associated amount of cryptocurrency to ensure a complete refund in case a conflict resolution ends in favour of the consumer (*collateral*).  
How to choose *v* depends on the identification system in use (see Section 4.3.1) and the resilience from providing, for free, a representative part of the dataset to one or multiple attackers. The analysis in Section 5.3 shows how to properly choose *v* based on the size of the dataset and the estimated amount of identities an attacker can hold.
3. **Provider**: Generate *s*.  
The provider should generate the symmetric encryption keys in such a way that in the case of disclosing some of them, the consumer will not be able to derive any other key (or the whole set). In addition, the consumer must be able to easily derive all the keys when they agree to buy the dataset. A simple way to achieve these features is by generating an initially private seed and to use a cryptographic hash function to compute the whole set of keys. For that reason, the seed *s* is calculated using a random number generator.
4. **Provider**: Compute  $K = [k_0 \dots k_{n-1}] \mid k_i = h(s + i) \forall i \in \{0 \dots n-1\}$ .  
The provider has to compute a set of *n* symmetric encryption keys ( $K = [k_0 \dots k_{n-1}]$ ). In DEFS, we compute each key as the hash function, for instance Keccak256, of the sum of the seed *s* and the index *i* using the following formula:  
$$k_i = \text{hash}(s + i) \forall i \in \{0 \dots n-1\}.$$
  
This construction has the expected properties: without the seed *s* the consumer cannot derive any other key, but once the seed is known, it is easy for the consumer to calculate the whole sequence of keys.
5. **Provider**: Generate *MHT(K)*.  
The provider builds a binary MHT for the set of keys, which is going to be used to generate the proof of the correctness of the keys used to encrypt the data portions. We denote the MHT of encryption keys as *MHT(K)*, its root *MRK*, and we refer to a membership proof of a leaf *i* as *MPK<sub>i</sub>*. Figure 4 shows an example of a *MHT(K)*. An exhaustive explanation about the algorithm to construct these trees can be found in [25]. Here, we will simply include a brief summary of this algorithm. To construct the *MHT(K)*, keys must be sorted by using their indexes, from 0 to *n* − 1. The leaves of the tree are calculated by hashing the keys in their respective position (hashed keys). The rest of the intermediate nodes in upper levels are just calculated by hashing the concatenation of the lower left and right nodes of the same branch. The tree construction continues until reaching the top level, in which we obtain the *MRK*. Notice that the *MRK* is the digest of the complete key set *K*, and it can be used as a proof of its correctness. It is also important to remark that the algorithm to construct this tree should be public, and all the entities have to use the same algorithm because any change in the keys or in the order of constructing it will cause an avalanche effect that will result in the root *MRK* being completely different.
6. **Provider**: Generate  $D = [d_0 \dots d_{n-1}] \mid d_i = \text{concat}(i, \text{data}_i) \forall i \in \{0 \dots n-1\}$ .  
Now, using the pre-existing data to be exchanged, the provider has to build an array ( $D = [d_0 \dots d_{n-1}]$ ) with the portions where each *d<sub>i</sub>* has the corresponding data and the index as a header  $d_i = \text{concat}(i, \text{data}_i)$ . This format is going to allow the

smart contract to determine if the cryptograms  $d_i$  have been properly generated. It is also important to note that each data portion  $data_i$  contains a group of random registries, not consecutive ones. To do that, registries are sorted using an external random generator tool provided by the marketplace, which should be open source and auditable to avoid duplicated entries in the same portion. This will avoid a potential attack where the consumer replicates several data requests to obtain free samples without committing any payment. More details about potential attacks can be found in Section 5.3.

7. **Provider:** Encrypt data.  $C = [c_0 \dots c_{n-1}] \mid c_i = E_{k_i}(d_i) \forall i \in \{0 \dots n-1\}$ .  
Now the provider is able to encrypt all the portions of the dataset  $d_i$ , and obtain the set of cryptograms  $C$ .
8. **Provider:** Generate  $MHT(C)$ .  
The provider builds a binary MHT for the set of cryptograms, which is going to be the proof of their correctness. We denote the MHT of cryptograms as  $MHT(C)$ , its root  $MRC$ , and we refer to a membership proof of a leaf  $i$  as  $MPC_i$ . Figure 5 shows an example of an  $MHT(C)$ . The algorithm of the  $MHT(C)$  is the same as the  $MHT(K)$ , but just changing the information used to construct it. Cryptograms are sorted by using their indexes, from 0 to  $n-1$ . The leaves of the tree are calculated by hashing the cryptograms in their respective position (hashed encrypted data). The rest of the intermediate nodes in upper levels are just calculated by hashing the concatenation of the lower left and right nodes of the same branch. The tree construction continues until reaching the top level, in which we obtain the root of the Merkle hash tree of cryptograms ( $MRC$ ). Note that the  $MRC$  is the digest of the complete set  $C$ , and it can be used as a proof of its correctness.
9. **Provider**→**SC:** Deploy smart contract with parameters ( $id$ ,  $n$ ,  $p$ ,  $v$ ,  $MRC$ ,  $MRK$ ,  $collateral$ ).  
Next, the provider deploys a smart contract in a ledger that stores the data-exchange identifier ( $id$ ), the total number of portions ( $n$ ), the number of sample portions to be revealed before the payment, the price ( $p$ ) of the dataset, and the amount of cryptocurrency to assure the fairness from the provider on the conflict resolution process ( $collateral$ ), and the root of both Merkle hash trees ( $MRC, MRK$ ). These roots will allow proving whether an element is or is not a cryptogram or a key and its position in the MHT. Using this, the system is able to efficiently assure the consumer that the provider cannot alter the committed dataset.
10. **Provider**→**Consumer:** Send  $id$  and the complete set of cryptograms ( $C$ ).  
Finally, the provider delivers the smart contract address, the data-exchange identifier  $id$ , and the complete set of cryptograms  $C$  to the consumer. Obviously, it would be totally impractical to exchange that amount of data using the ledger as storage. Instead, the exchange of cryptograms between the provider and consumer is done off-blockchain. Notice also that since the seed will be public at the end of the process, all the off-blockchain traffic must have been exchanged using a secure channel.
11. **Consumer**→**SC:** Request data from  $id$ .  
At this point, the consumer has the  $id$  which is related to the deployed smart contract and can read the values set by the provider in step 9.
12. **SC**→**Consumer:** Reply with values  $n$ ,  $p$ ,  $v$ ,  $MRC$ ,  $MRK$ .  
Now the consumer has the total number of portions ( $n$ ), the number of samples they can obtain before committing the payment ( $v$ ), the price of the data ( $p$ ), the Merkle roots of both trees ( $MRC, MRK$ ), and the complete set of cryptograms ( $C$ ).
13. **Consumer:** Compute  $MHT(C)$  and verify  $MRC$ .  
As the consumer has the complete set of cryptograms  $C$ , he has the capability and responsibility to re-generate  $MHT(C)$  to verify that the root  $MRC$  calculated is coherent with the one at the smart contract. If so, the consumer knows that all cryptograms  $c_i \forall i \in \{0 \dots n-1\}$  were properly generated and match the  $MRC$ . The consumer is responsible for verifying the  $MHT(C)$  at this very moment, and if he continues with

the protocol, tacitly accepts the correctness of the generation of the cryptograms. This means that in the case of later conflict resolution, the consumer cannot argue that the cryptograms were wrongly generated to get a refund.

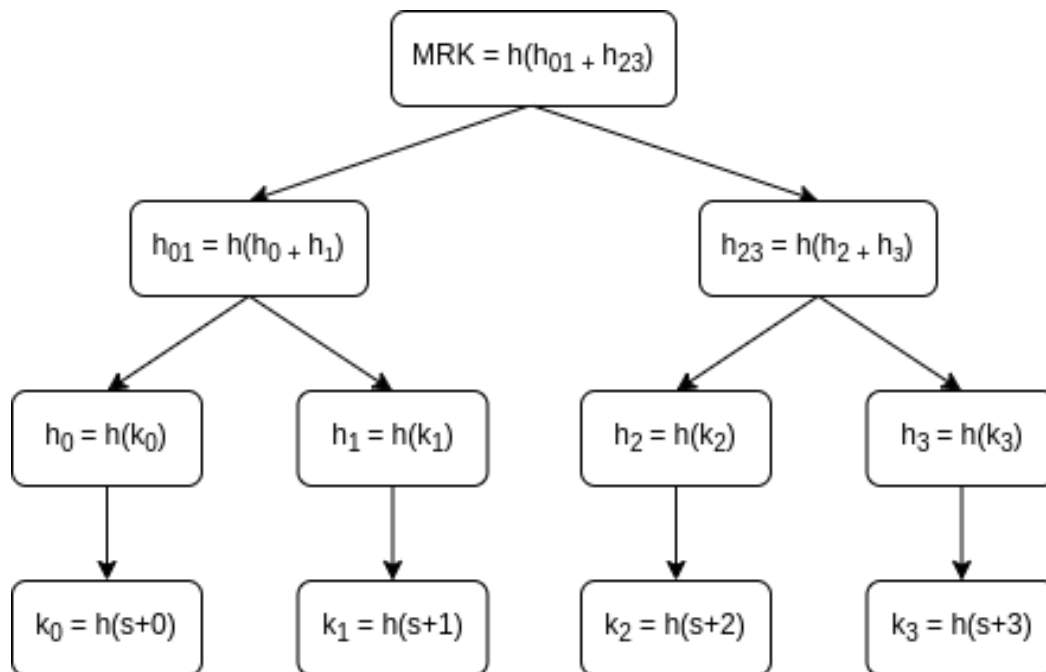


Figure 4. Tree of keys  $MHT(K)$ . In this example,  $MRK = H_{0123}$ .

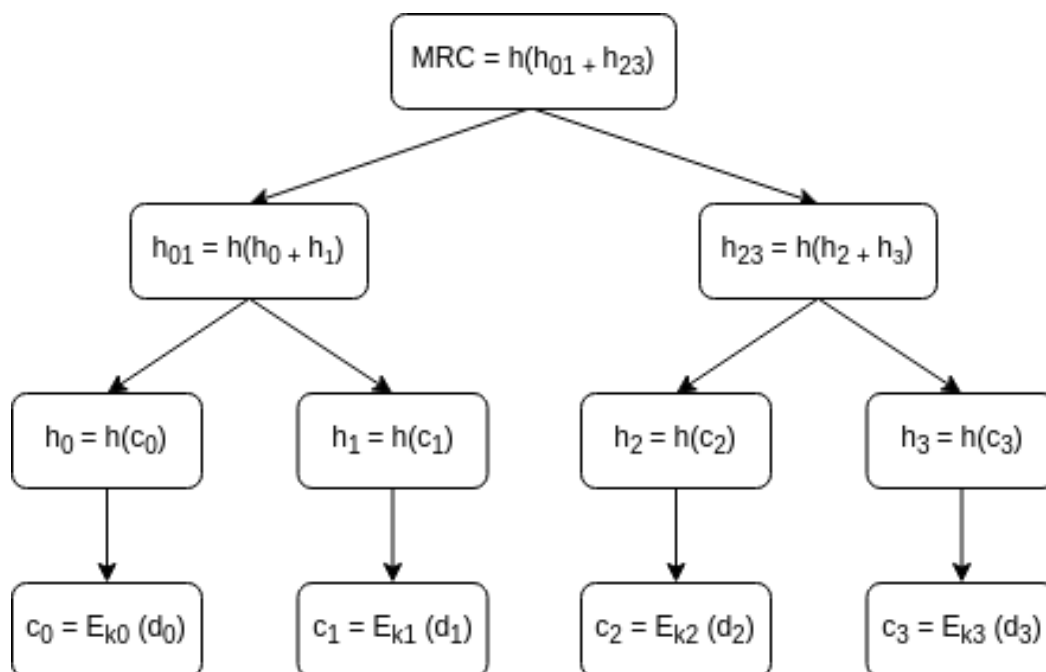


Figure 5. Tree of cryptograms  $MHT(C)$ . In this example,  $MRC = H_{0123}$ .

#### 4.3.4. Protocol Execution

Once the protocol preparation phase is completed, the consumer requests the provider to reveal some sample portions. If these samples convince the consumer about the quality of the dataset, the consumer commits the payment. It is important to note that if this

protocol execution phase ends as expected, the protocol is completed, and there is no need to execute the conflict resolution procedure.

Figure 6 details the interactions between provider, consumer, and smart contract during the execution phase of the DEFS protocol. The steps of this phase are enumerated as:

1. **Consumer:** Generate  $R = [r_0 \dots r_{v-1}] \mid \forall r_i \in \{0 \dots n-1\}$ .  
The consumer selects, at random, the set of indexes  $R$  which correspond to the sample portions to be revealed (for free). We consider that the provider should not decide on its own which data samples will be revealed because they could decide to use a biased (not fair) set of samples. The consumer is also not able to purposely choose particular registries because the dataset was previously randomly sorted by the provider. Therefore, the consumer will choose an array of  $v$  values, at random, within the range 0 to  $n-1$ , corresponding to the indexes of the sample portions. As none of the entities control which registries are going to be disclosed, the fairness of this process is assured.
2. **Consumer**→**Provider:** Request  $R$  which contains  $v$  indexes to be revealed.  
The consumer informs the provider of the  $v$  indexes of the sample portions to be revealed ( $R$ ).
3. **Provider**→**Consumer:** Return  $k_i$  and  $MPK_i \forall i \in \{R\}$ .  
The provider discloses the keys associated with the  $v$  indexes and the Merkle proofs to verify them in the  $MHT(K)$ . This process is done totally off-blockchain. Note that allowing the consumer to choose  $v$  sample portions with no cost could result in attacks as the consumer could try to get a large amount of free data by repeating the process of getting small samples. We do not consider this attack especially dangerous because the consumer cannot choose particular registries of the dataset. In addition, the provider can decide how many times consumers are allowed to get a sample set without making a final deal, and they can use the marketplace to blacklist abusive consumers. Nonetheless, a comprehensive analysis of potential attacks is made in Section 5.3. As shown in the analysis, the provider should carefully choose  $v$  and  $n$  to minimise the impact of such attacks.
4. **Consumer:** Verify each proof.  $MPK_i \forall i \in \{R\}$ .  
The customer should verify that the proofs sent by the provider match the  $MRK$ . To provide a partial example, consider the case of Figure 4, in which  $n = 4$ . Let us consider the case of one single sample portion  $v = 1$ , and the consumer has chosen index 3 to disclose. In this case,  $k_3$  is sent to the consumer, and the proofs of the correctness are  $MPK_3 = (h_{01}, h_2)$ . The consumer should verify if the following expression matches:

$$\text{hash}(\text{concat}(h_{01}, \text{hash}(\text{concat}(h_2, \text{hash}(k_3))))) == MRK$$

Note that the expression is just calculating the root of  $MHT(K)$ , and comparing it with the  $MRK$  value published in the smart contract. If the values match, the key can be considered valid as it matches the proofs, and the consumer will continue with the protocol. If not, that means that the provider did not send the proper proofs and that the key is not verifiable. In this case, the consumer can abandon the protocol.

5. **Consumer:** Decrypt and verify each data sample.  
Now that the consumer is sure that the  $v$  keys are valid, he can decrypt the sample portions with the received keys:

$$d_i = E_{k_i}^{-1}(c_i) \forall i \in R$$

Once this is done, the consumer has to verify that the resulting data portions have the expected format:

$$d_i = \text{concat}(i, \text{data}_i)$$

If not, the data samples are invalid, and the consumer can end the protocol at this moment. If the format is valid, the consumer will evaluate the data samples  $data_i \forall i \in R$ . If the samples do not convince the consumer to pay for the whole dataset, the protocol ends here, but if they do, they will continue with the following steps.

6. **Consumer**→**SC**: Transaction committing payment.  
If the samples convinced the consumer, he will send a transaction to the smart contract with the payment ( $p$ ) to buy the dataset.
7. **Consumer**→**SC**: Subscription to seed revelation.  
The consumer also subscribes to the 'seedReleased' event, expecting to receive a notification when the provider publishes the seed to allow complete data decryption.
8. **(Timeout 1) Consumer**→**SC**: Seed not released and the consumer is refunded.  
If the provider has not released the seed in time, a first timeout (Timeout 1) will expire, and after that, the smart contract will allow the consumer to refund the payment. This is an unhappy path in the protocol.
9. **Provider**→**SC**: Transaction publishing the seed.  
The provider discloses the seed value  $s$  via a blockchain transaction before Timeout 1 expires. This is the happy path of the protocol.
10. **SC**→**Consumer**: Event to consumer about seed revelation.  
As result of executing the transaction, the smart contract generates an event, and the consumer will be notified that the seed value  $s$  has been revealed. The smart contract stops Timeout 1 (due to seed revelation), and starts Timeout 2, allowing the consumer to start conflict resolution. Once that the consumer has the seed, they can derive all the keys:

$$k_i = \text{hash}(s + i) \forall i \in \{0..(n-1)\}$$

Now the consumer has all the cryptographic material to decrypt the cryptograms  $C$  but, previously, they had to verify that all the keys are correct. The procedure is similar to the one performed in Step 4, but with the difference that now the consumer has the capacity of re-generating the whole  $MHT(K)$ . If the consumer detects that (one or several) keys were not properly generated, they can start the optional phase of conflict resolution to obtain a refund. On the contrary, if all the keys were properly generated, the consumer can decrypt the previously received cryptograms  $C$  and access the whole dataset:

$$d_i = E_{k_i}^{-1}(c_i) \forall i \in \{0..(n-1)\}$$

The consumer also has to verify that the  $d_i$  have the proper format, in the same manner that was done in Step 5 but for all the cryptograms. If the consumer detects that one or more cryptograms were not properly generated (they do not have the proper format), they can start the optional phase of conflict resolution to obtain a refund. Note that at this point, the consumer cannot argue that the cryptograms do not match the  $MHT(C)$  because this should have been verified in the protocol preparation phase. On the contrary, if all the cryptograms (once decrypted) have the proper format, the consumer has the complete dataset and can consider the protocol ended.

11. **(Timeout 2) Provider**→**SC**: Withdraw and protocol end.  
If Timeout 2 expires, it means that the consumer considers that the keys were properly generated and the cryptograms have the proper format (because if not, they would have previously started the conflict resolution). In this case, the provider can send a transaction to the smart contract to withdraw the payment and end the protocol.

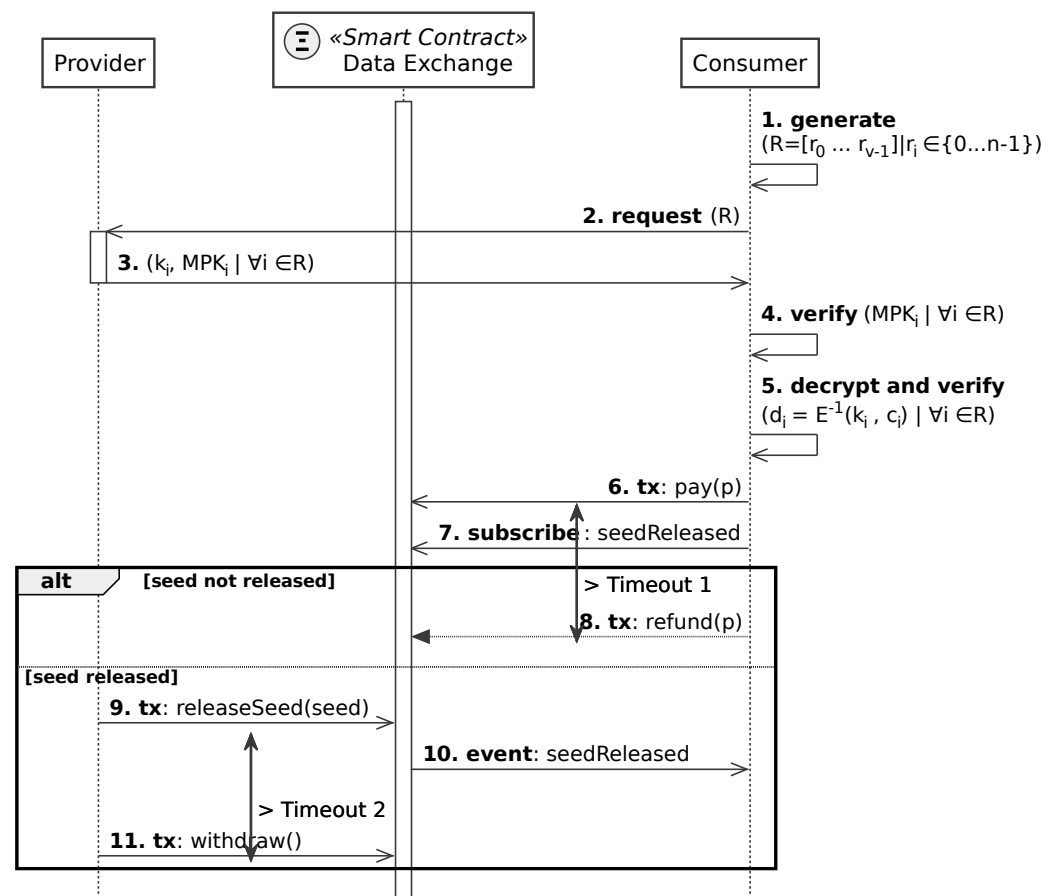


Figure 6. Protocol execution: sequence diagram.

#### 4.3.5. Conflict Resolution

The final aspect that our protocol has to solve is conflict resolution, which may appear if the consumer detects misbehaviour. As previously stated, the consumer must start the conflict resolution phase before the expiration of Timeout 2 in the protocol execution phase. If Timeout 2 expires, the smart contract considers the protocol ended and the provider can receive payment.

We will consider two cases that are relevant and can end with a refund if the consumer is able to demonstrate the misbehaviour: (1) a key is not properly generated; and (2) a decryption of a cryptogram does not have the proper format. There are also a couple of extra cases for dispute that will not end with a refund: (3) cryptograms are not properly generated; and (4) dataset is of bad quality.

##### 1. A key is not properly generated:

When consumers obtain the seed  $s$  in Step 9 of the protocol execution phase, they are able to generate the whole set of keys  $K$ . The way to check if the set of keys is compliant or not is by generating the whole set  $K$  with the formula  $k_i = \text{hash}(s + i)$  and also re-constructing the whole  $MHT(K)$  and verifying that the calculated root matches the one published in the smart contract. At this moment, a consumer can know that the registered MRK is incorrect. However, in general, he cannot detect which keys were not properly generated. Although not possible in general, there are particular cases in which the consumer can do the detection of wrong keys. The detection is possible when the wrong keys are either one of the keys used for encrypting the samples, or a sibling key of them. In both cases, the consumer has got an MPK from the provider that matches the registered MRK. Then, if one of those keys does not follow the agreed format  $\text{hash}(i + s)$ , the consumer can send the hash of the wrong key, its MPK, and the index in conflict to prove to the smart contract that



the provider committed an incorrect MRK. For simplicity, we will assume that the consumer detects one single not-compliant key  $k_e$ , but this discussion is completely valid in the case of having multiple not compliant keys (DEFS does not distinguish between one or several not compliant keys, and in the case where it is demonstrated that one single key is not properly generated, the entire payment will be refunded to the consumer). The following are the steps that the consumer, smart contract, and provider have to follow during conflict resolution about a specific key  $k_e$ . The associated sequence diagram is detailed in Figure 7:

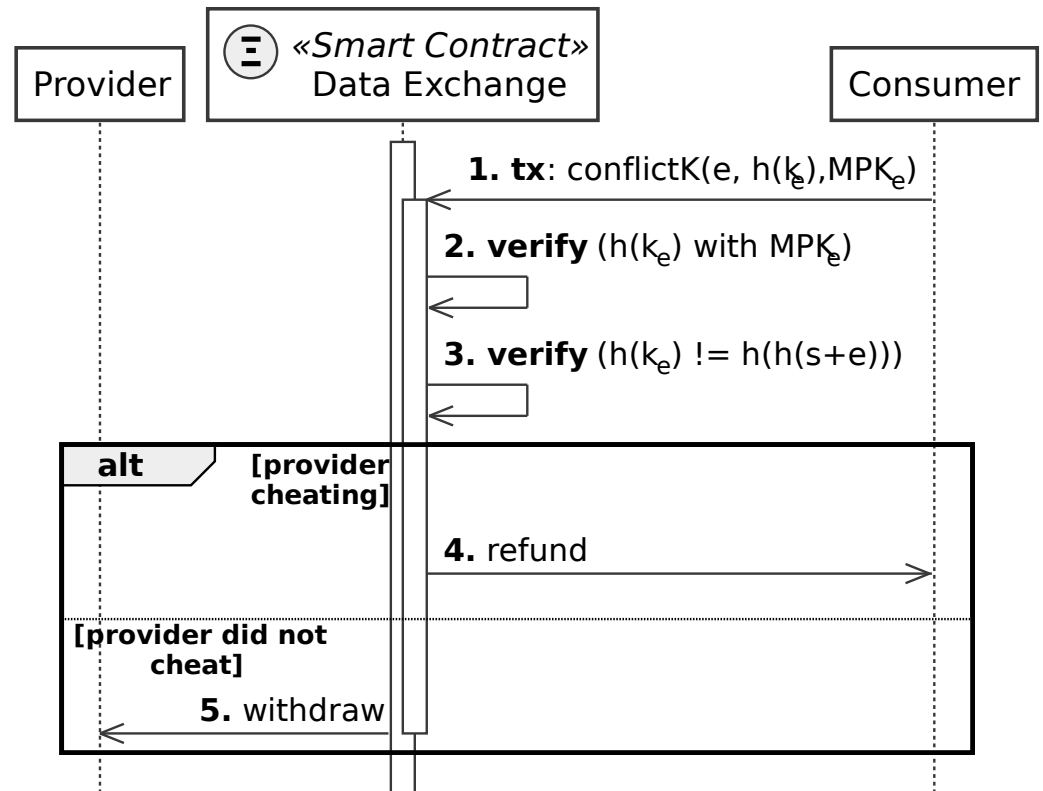


Figure 7. Protocol resolution-k: sequence diagram.

- (1) **Consumer**→**SC**: tx: conflictK( $e, h(k_e), MPK_e$ )  
This conflict resolution is performed with a transaction from the consumer to the smart contract calling the function 'conflictK'. In this transaction, the consumer sends as parameters the problematic index ( $e$ ), the hash of the invalid key ( $h(k_e)$ ), and its associated proof ( $MPK_e$ ).
- (2) **SC**: Verify  $h(k_e)$  with the  $MPK_e$   
The smart contract verifies that  $h(k_e)$  and  $MPK_e$  match the MRK from Step 9.
- (3) **SC**: Verify  $h(k_e) \neq h(h(s+e))$   
The smart contract verifies whether  $h(k_e)$  matches  $h(h(s+e))$  or not. If there is a match, it means that the provider did not cheat, while if the check does not match, then it means that the provider cheated.
- (4) **SC**: Provider cheating  
If the provider cheated, the smart contract refunds the consumer. In this case, the consumer receives the price of the data ( $p$ ) and also the cost of the transactions he sent. The cost of the transactions is taken from the provider's collateral.
- (5) **SC**: Provider not cheating  
If the provider did not cheat, the smart contract automatically transfers the payment ( $p$ ) to the provider.

## 2. Cryptograms do not have the proper format:

This situation happens when there is a conflict in  $D$ , and so one or several data portions do not have the proper format. Just a remark that there was a previous checking of this type in Step 5 of the protocol execution phase, in which  $v$  of the possible sample portions were tested to see if they had the proper format:

$$d_i = \text{concat}(i, \text{data}_i)$$

However, not all the data portions were tested (only  $v$  of  $n$ ). For simplicity, we will assume that there is one single not-compliant portion  $d_e$ , but this discussion is completely valid in case of having multiple not compliant portions.

This scenario implies that the decryption of a  $c_e$  results in a  $d_e$  that does not correspond with the expected format  $d_e = \text{concat}(e, \text{data}_e)$ . Specifically, the decrypted cryptogram does not start with the expected index ( $e$ ). In this case, the consumer can start the conflict resolution about the format of data. Figure 8 shows the sequence diagram about this scenario.

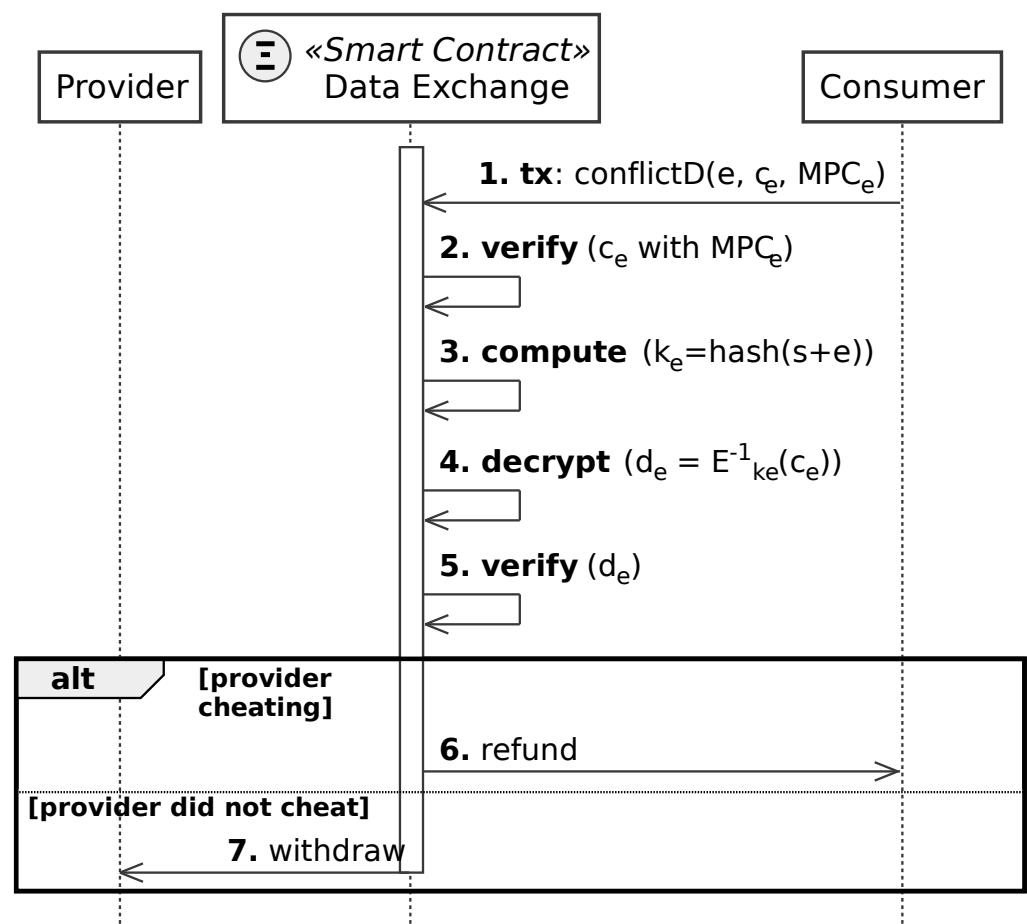


Figure 8. Protocol resolution-d: sequence diagram.

### (1) tx: conflictD:

The conflict resolution about a  $c_e$  starts with a transaction from the consumer to the smart contract. To do that, the transaction calls the smart contract function 'conflictD' and sends as parameters the problematic index ( $e$ ), the cryptogram involved  $c_e$ , and the proofs of the validity of this particular cryptogram. The intention of the consumer is to show that the cryptogram  $c_e$  was properly generated by the provider, but that after decrypting it with  $k_e$ , the resulting data portion has a bad format.

In this case, the protocol can be resolved in a single transaction from the consumer, because the smart contract can compute  $k_e$ . Note that there is no need for an extra timeout because the provider does not need to send anything, and all the proofs are available for the smart contract to compute.

(2) **Verify  $MPC_e$ :**

The smart contract has to verify that the cryptogram provided by the consumer is valid, that is to say, that calculating  $MRC$  for the problematic cryptogram  $c_e$  from  $MPC_e$  is coherent with the root stored in the smart contract. To provide some piece of example, consider the case of Figure 5, in which  $n = 4$ . Let us consider the case of demonstrating that  $c_3$  is properly generated. In this case, the proofs of the correctness of  $c_3$  are  $MPK_3 = (h_{01}, h_2)$ . The smart contract should verify if this expression matches:

$$\text{hash}(\text{concat}(h_{01}, \text{hash}(\text{concat}(h_2, \text{hash}(c_3))))) == MRC$$

Notice that the expression is just calculating the root of  $MHT(C)$ , and comparing it with the  $MRC$  value published in the smart contract. If these values match, the cryptogram can be considered valid as it matches the proofs, and the protocol continues with the following step. If not, the consumer did not send the proper proofs to demonstrate that the provider was cheating and that the withdrawal of the money is automatically received as shown in Step (7).

(3) **Compute  $k_e$ :**

Now the smart contract knows that the cryptogram  $c_e$  is valid. Next, the smart contract computes the associated key  $k_e = \text{hash}(s + e)$ . Remember that  $s$  was published by means of the transaction sent in Step 9 of the protocol execution phase.

(4) **Decrypt  $c_e$ :**

The smart contract has the valid key and the valid cryptogram, so it is able to decrypt and obtain  $d_e = E_{k_e}^{-1}(c_e)$ .

(5) **Verify  $d_e$  format:**

The smart contract can verify if the data portion  $d_e$  has the correct format.

(6) **refund:**

If the data portion  $d_e$  does not start with the index  $e$ , the provider was cheating, so the transaction ends transferring the costs ( $p + \text{collateral}$ ) to the consumer.

(7) **withdraw:**

If the data portion  $d_e$  start with the index  $e$ , the provider was not cheating, and the transaction ends transferring the payment ( $p$ ) and the *collateral* to the provider.

3. **Cryptograms are not properly generated:**

This case happens when there is a conflict in  $C$  and one or several cryptograms do not match the root  $MRC$  of the tree  $MHT(C)$ . As previously stated, the consumer receives all the cryptograms  $C$  in Step 10 of the protocol preparation and verifies the correctness of the whole set of cryptograms in Step 13. The consumer was responsible for verifying the  $MHT(C)$  at this very moment, and if any problem during this checking was found, the protocol will simply be aborted before committing any payment. However, if the consumer continued with the protocol, they were tacitly accepting the correctness of the generation of the cryptograms and the corresponding  $MHT(C)$  and  $MRC$ . In case the consumer detects a cryptogram ( $c_e$ ) not matching the  $MRC$  at the protocol execution or protocol resolution phases, they cannot try to get a refund, and for this reason, the conflict resolution is not considering that case.

4. **Dataset is of bad quality:**

This case happens when the consumer obtains a valid data portion  $d_e$  (with the correct format  $d_e = \text{concat}(e, \text{data}_e)$ ), but the content does not have the quality expected by the consumer. In this particular case, the DEFS protocol is not able to consider the quality

of the dataset (most probably, assessing the goodness of a dataset requires human interaction and cannot be made automatically by the smart contract), so this case is out of its scope of discussion and no refund can be requested from the consumer's side. The  $v$  sample portions that were disclosed in Step 5 of the protocol execution (for free, prior to any payment) alleviate this possibility. In any case, the marketplace can consider having a reputation tool to value data providers and try to avoid this kind of behaviour.

#### 4.3.6. State Diagram

The protocol operation and the interactions between the different stakeholders and the smart contract are detailed in Figure 9.

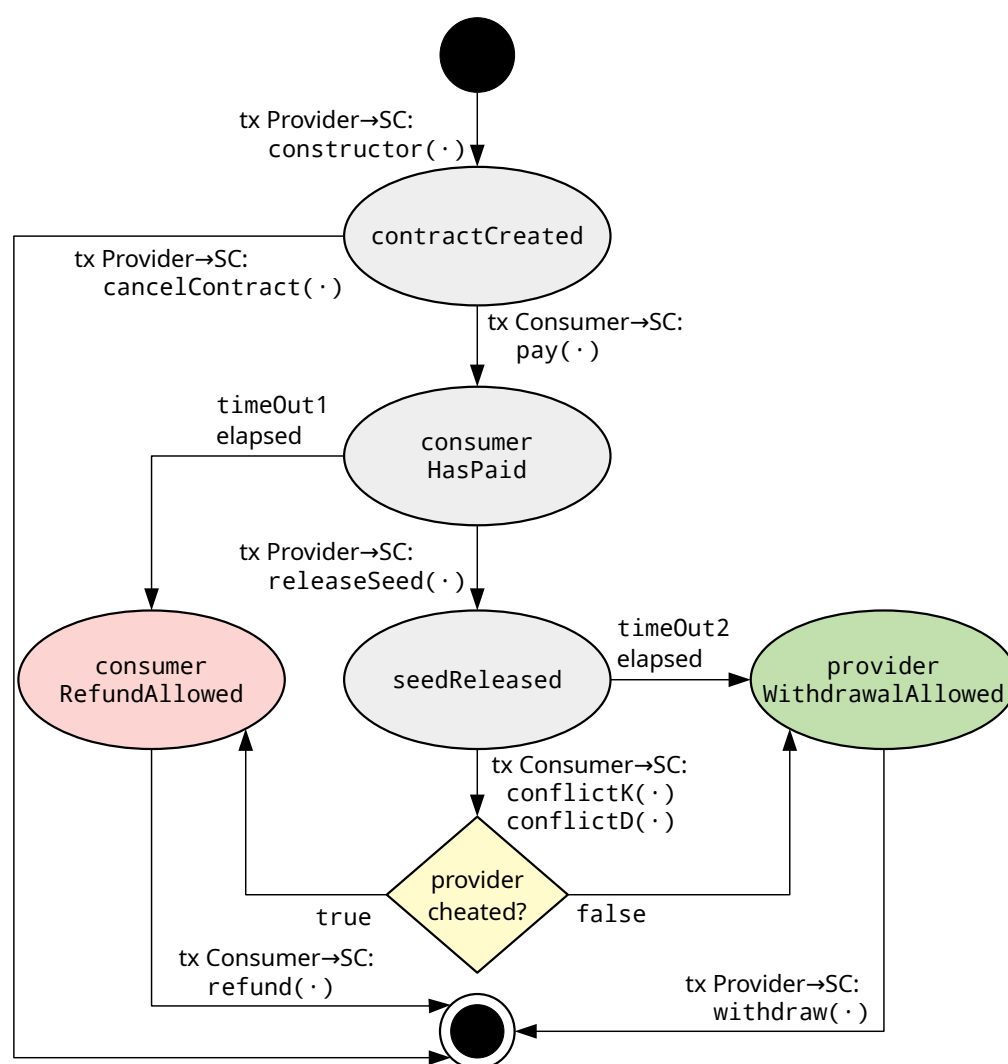


Figure 9. State diagram of the smart contract.

#### 5. Security Analysis

The DEFS enables providers and consumers to make commercial agreements via a smart contract. Essentially, the smart contract ensures that if a consumer agrees to purchase a dataset based on the provided random samples, then the provider receives the right amount of money and the consumer gets access to the whole dataset.

### 5.1. Protection against Channel Attacks

Channel attacks are those in which the attacker purposely tries to eavesdrop information from the channel. In our proposal, two different channels must be considered: the on-chain channel with the interactions with the smart contract, and the off-chain channel between the data provider and the data consumer.

It is assumed in this work that the off-chain channel is protected from overhearing and tampering, for instance, by forcing it to use transport layer security (TLS). As a result, no information would be exposed in this channel. However, the on-chain channel is public and all the interactions with the smart contract will be available for an attacker, namely the data-exchange identifier  $id$ , number of data portions  $n$ , the price of the dataset  $p$ , and the roots of the MHT of cryptograms  $MRC$  and keys  $MRK$ .

The only valuable information an attacker could obtain is the size  $n$  and price  $p$  of the dataset. The data-exchange identifier  $id$  is just an identifier. The  $MRC$  and  $MRK$  leak no information regarding the keys and cryptograms since the cryptographic hash function used to create the MHT is assumed to be preimage- and collision-resistant.

### 5.2. Consumer's Protection against Provider Attacks

In the first part of the protocol, the consumer gets access to a set of free random samples from the dataset. The provider commits the structure of cryptograms and key structure (when committing the  $MRC$  and  $MRK$  to the smart contract). Then, the consumer can request a specific set of random samples before making an assessment with regard to buying the dataset or not.

A malicious provider would like to send a selected set of samples that make the dataset more appealing. However, the provider has no control over the selected samples and changing them on the fly would require that the fake sample keys and cryptograms collude with the committed ones. This attack is assumed to be unfeasible since the probability of finding collisions in the cryptographic hash function used to generate the Merkle trees is assumed to be negligible.

Providers would also hold datasets with both good and bad data, meaning bad fake or even duplicated data. However, since providers cannot choose the requested free samples, bad samples could be detected during the evaluation of the free samples. In any case, it is up to the consumer to decide if the amount of samples is representative enough to get a fair idea of the content.

Finally, another potential attack would be that of a provider releasing wrong keys, wrong cryptograms, or incorrectly encrypted data after getting paid. In these cases, the consumer can make use of the different conflict resolutions explained in Section 4.3.5 and get a refund. Recall that, in the case of incorrect cryptograms, it is important that the consumer validates all cryptograms before doing the payment, since otherwise he will not have access to the conflict resolution.

### 5.3. Provider's Protection against Consumer Attacks

Every time a consumer engages in the protocol, he receives  $v$  samples of the product. A dishonest consumer with several identities could accumulate free samples and try to get the whole dataset without paying for it or, equivalently, collaborate with other consumers to get as many data samples as possible for free. As we show hereunder, the provider can adjust the amount of samples disclosed to the consumer to reduce the probability of these attacks succeeding.

Before analysing a general setup, let us consider a simple example in which the provider has 2 different samples and every consumer gets 1 for free. That is,  $n = 2$  and  $v = 1$ . Let us compute what is the probability that a malicious consumer with several identities gets both samples. With one interaction, the consumer gets only half of the product for free. If the consumer creates a new identity, he would get a new sample, but he would only get the whole product if the new sample is different from the previous one. Hence, we need to compute what is the probability that a new sample is different

from the previous one. We can think of this scenario as tossing two coins and finding the probability that we get a head ( $H$ ) and a tail ( $T$ ), which is  $1/2$ . Therefore, with two identities, a dishonest consumer would have a 50% chance to get the whole set of samples. If the consumer creates a new identity and gets a third sample, the probability of them getting the two different samples goes up to 75%, which is the probability that in a sequence of 3 coin tosses, at least one is a head and another a tail.

In general, it is assumed that a consumer gets  $v$  random samples and is able to interact with the provider with  $k$  different identities, obtaining a total of  $m = k \times v$  free random samples. Like before, the probability that among  $m$  samples,  $n$  of them are different, is the same as the probability of having  $n$  different elements in a sequence of  $m$  elements. If  $m < n$ , then clearly this probability is 0. If  $m = n$ , then there are  $n!$  different sequences with the  $n$  elements. Hence, the probability of getting a sequence of this kind is  $n!/n^n$ . When  $n$  is large enough, this probability is very low. To get an idea, we show in Table 2 what happens if the provider discloses 10% of samples of his dataset and the consumer creates 10 identities to get a total of  $n$  samples. Note that, even for small values of  $n$ , the probability that a consumer gets the whole dataset is very low.

**Table 2.** Probability that the consumer gets the whole dataset if the provider makes 10 interactions with the consumer and discloses 10% of his dataset every time.

$n$	$v$	$k = m/v$	Probability of Getting the Whole Dataset
10	1	10	$3.62 \times 10^{-4}$
100	10	10	$9.33 \times 10^{-43}$
1000	100	10	$4.02 \times 10^{-433}$
10,000	1000	10	$\sim 10^{-4340}$
100,000	10,000	10	$\sim 10^{-43426}$

To increase the odds of getting all samples, a dishonest consumer would create more identities so that  $m \geq n$ . In this scenario, we need to calculate what is the probability that among the  $n^m$  possible outcomes, the consumer gets  $n$  distinct data samples. If we go back to the case  $n = 2$ ,  $v = 1$ , and a consumer with three identities ( $k = 3$ ), we should count how many sequences of three elements contain 2 distinct elements. Or equivalently, what is the probability that after three coin tosses, we get at least one head and one tail. We can think of this problem as counting the different ways in which we can assign the three positions of the sequence  $\{1, 2, 3\}$  to a head or a tail such that at least one is a head and another a tail. This counting is precisely the number of ways in which we can partition the set  $\{1, 2, 3\}$  into two non-empty sets:  $\{1, 2\} \cup \{3\}$ ,  $\{1, 3\} \cup \{2\}$ , and  $\{2, 3\} \cup \{1\}$ . If we assign the first set to heads and the second set to tails, the partitions lead to the three sequences  $HHT, HTH, TTH$ , and if we do the opposite assignment, we get  $TTH, THT, HHT$ . As a result, we get a total of  $3 \times 2 = 6$  different sequences containing at least one head and one tail, which divided by the  $2^3 = 8$  possibilities, results in the 75% chances we claimed before.

In general, the number of sequences of  $m$  elements that contain  $n$  distinct elements is equivalent to the number of ways we can partition the  $m$  positions of the sequence into  $n$  non-empty sets multiplied by the number of permutations of  $n$  distinct elements. This count is precisely the Stirling number of the second kind  $S(m, n)$  multiplied by the number of permutations of  $n$  elements [26]:

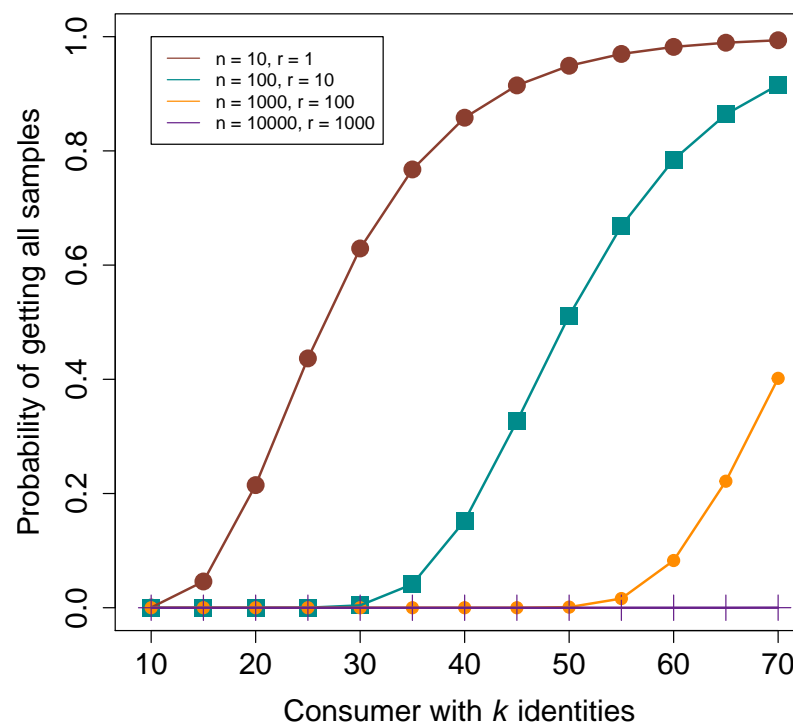
$$S(m, n)n! = \sum_{j=0}^n \binom{n}{j} (-1)^{n-j} j^m.$$



Hence, the probability that a consumer with  $k = m/v$  identities gets the whole dataset is

$$\frac{S(m,n)n!}{n^m} = \frac{1}{n^m} \sum_{j=0}^n \binom{n}{j} (-1)^{n-j} j^m.$$

In Figure 10, we illustrate the probability that a consumer with  $k$  identities that gets 10% of free samples with each identity, obtains the whole dataset of  $n$  registries.



**Figure 10.** Probability that a consumer with  $k$  different identities that gets a 10% of samples with each identity obtains all  $n$  different data samples.

Note that a provider with a dataset of 10,000 registries disclosing a random 10% of it for free should almost only be worried about consumers that are able to create more than 50 different identities. The off-chain identification system (see Section 4.3.1) should be chosen to minimise or even impede the likelihood of an attacker getting more than  $k$  identities. Typically,  $n$  is several orders of magnitude higher than 10,000, so the amount of fake identities needed to perform this attack would make it infeasible in practice.

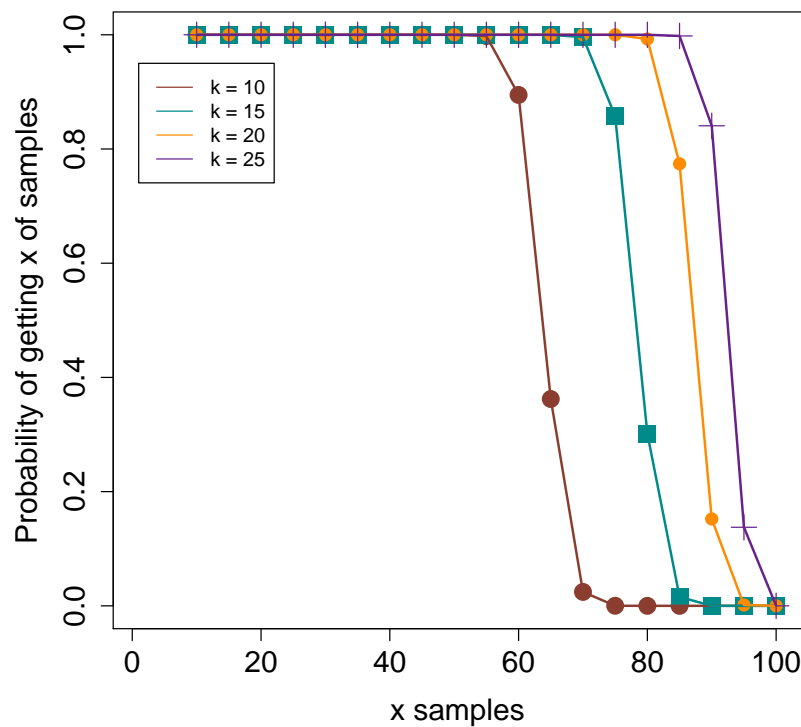
Even if the probability of getting the whole dataset is low enough, the provider may also want to avoid disclosing a large valuable set of data. We analyse the probability of obtaining a meaningful amount of samples for different values of  $m$ . In the same way as we argued before, the probability of getting  $x$  different samples is the number of combinations of sequences of  $m$  elements with  $x$  distinct elements, but now multiplied by the combinations of  $x$  elements that we can make with  $n$  distinct elements. That is,

$$P(\text{get exactly } x \text{ distinct samples}) = \frac{S(m,x)x!}{n^m} \binom{n}{x} = \frac{1}{n^m} \binom{n}{x} \sum_{j=0}^x \binom{x}{j} (-1)^{x-j} j^m.$$

Therefore, the probability of getting  $x$  distinct elements in a sequence of  $m$  elements is

$$\begin{aligned} P(\text{get at least } x \text{ distinct samples}) &= 1 - \sum_{i=0}^{x-1} P(\text{get exactly } i \text{ distinct samples}) \\ &= 1 - \sum_{i=0}^{x-1} \binom{n}{i} \frac{S(m, i) i!}{n^m}. \end{aligned}$$

To illustrate the tendency of these probabilities, we have depicted in Figure 11 the probability of obtaining at least 5, 10, 15, ..., 100 different samples from a dataset with 100 different samples. The different lines correspond to the probabilities using a different number of identities.



**Figure 11.** In a database of  $n = 100$  portions and  $v = 10$  free samples per consumer, this graphic depicts the probability of getting a different portion of samples for free with several identities:  $k = 10$  (large red circles),  $k = 15$  (blue squares),  $k = 20$  (small orange circles),  $k = 25$  (purple lines).

As we can see, there is a significant drop in the probability of obtaining at least a 65% of samples with 10 identities, but with more identities, this inflection point moves up to 80% (with  $k = 15$ ), 90% (with  $k = 20$ ), and 95% (with  $k = 25$ ). Therefore, even though the probability of obtaining the whole 100% of the dataset is very close to 0, with 25 identities, it is possible to obtain at least the 90% of it with probability 0.84.

As we have shown, the provider can mitigate the risks of identity-replication by strongly authenticating consumers through the off-chain channel (see Section 4.3.1) and by adjusting the amount of free samples disclosed to the consumer.

In general, the provider should find a trade-off between securing the dataset while, at the same time, letting the consumer get a fair idea its content.

Lastly, we would like to remark that the provider is protected against a consumer that does not pay after acknowledging the protocol since the result is revealed only once the payment is in the smart contract, and the cryptograms are useless without it.

## 6. Conclusions

Distrust is one of the main obstacles to implementing exchanges between data providers and data consumers in a decentralised way. In this article, we present a protocol that allows a consumer to probabilistically obtain and check a subset of a dataset on sale from a provider before committing to payment. The protocol is executed using a smart contract deployed in a public distributed ledger. Once the consumer agrees to buy the dataset, the payment process, the agreed terms, and the possible refunds are managed and enforced by the smart contract. To expose the dataset, our protocol splits the data in portions and encrypts and stores each portion off-chain. Then, we created a MHT for the cryptograms and another MHT for the encryption keys. The encryption keys are related to each other using a cryptographic hash function in a way that allows us to implement a cost-efficient conflict resolution mechanism. The security analysis of our protocol shows that consumers and providers are economically protected and that the provider can reduce the risks of identity-replication attacks by adjusting the amount of free samples disclosed to the consumer.

**Author Contributions:** Conceptualisation, J.L.M.-T. and J.H.-S.; methodology, R.G.-D.; validation, R.G.-D., J.H.-S., O.E., M.B.-M. and J.L.M.-T.; formal analysis, M.B.-M.; investigation, R.G.-D. and J.L.M.-T.; data curation, O.E. and J.L.M.-T.; writing—original draft preparation, R.G.-D. and M.B.-M.; writing—review and editing, J.H.-S., O.E. and J.L.M.-T.; visualisation, R.G.-D. and M.B.-M.; supervision, J.L.M.-T.; project administration, J.H.-S., O.E. and J.L.M.-T.; funding acquisition, J.H.-S., O.E. and J.L.M.-T.; All authors have read and agreed to the published version of the manuscript.

**Funding:** This research has been funded by i3Market (H2020-ICT-2019-2 grant number 871754). This work is also supported by the TCO-RISEBLOCK (PID2019-110224RB-I00), ARPASAT (TEC2015-70197-R) and by the Generalitat de Catalunya grant 2014-SGR-1504.

**Conflicts of Interest:** The authors declare no conflicts of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## Abbreviations

The following abbreviations are used in this manuscript:

API	application programming interface
DEFS	data exchange with free sample
DLT	distributed ledger technology
IPFS	interplanetary file system
MHT	Merkle hash tree
MR	Merkle hash root
MP	Merkle proof
MPC	Merkle proof of cryptograms
MPK	Merkle proof of encryption keys
MRK	root of the Merkle hash tree of keys
MRC	root of the Merkle hash tree of cryptograms
TTP	trusted third party
SC	smart contract
SDTE	secure data trading ecosystem

## References

1. Gopalkrishnan, V.; Steier, D.; Lewis, H.; Guszcz, J. Big Data, Big Business: Bridging the Gap. In Proceedings of the 1st International Workshop on Big Data, Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models and Applications, Beijing, China, 12 August 2012; Association for Computing Machinery: New York, NY, USA, 2012; BigMine '12, pp. 7–11. [\[CrossRef\]](#)
2. Thomas, L.D.W.; Leiponen, A. Big data commercialization. *IEEE Eng. Manag. Rev.* **2016**, *44*, 74–90. [\[CrossRef\]](#)
3. Ravi, N.; Sunitha, N.R. Introduction of Blockchain to Mitigate The Trusted Third Party Auditing for Cloud Security: An Overview. In Proceedings of the 2017 2nd International Conference On Emerging Computation and Information Technologies (ICECIT), Tumakuru, India, 15–16 December 2017; pp. 1–6. [\[CrossRef\]](#)
4. Bashir, I. *Mastering Blockchain*; Packt Publishing Ltd.: Birmingham, UK, 2017.

5. Yaga, D.; Mell, P.; Roby, N.; Scarfone, K. Blockchain Technology Overview. 2019. Available online: <https://arxiv.org/abs/1906.11078> (accessed on 10 May 2021).
6. Wood, G. Ethereum: A secure decentralised generalised transaction ledger. *Ethereum Proj. White Pap.* **2014**, *151*, 1–32.
7. Haider, F. Compact Sparse Merkle Trees. Cryptology ePrint Archive, Report 2018/955. 2018. Available online: <https://eprint.iacr.org/2018/955> (accessed on 1 June 2021).
8. Dahlberg, R.; Pulls, T.; Peeters, R. Efficient Sparse Merkle Trees: Caching Strategies and Secure (Non-)Membership Proofs. Cryptology ePrint Archive, Report 2016/683. 2016. Available online: <https://eprint.iacr.org/2016/683> (accessed on 1 June 2021).
9. Al-Kuwari, S.; Davenport, J.H.; Bradford, R.J. Cryptographic Hash Functions: Recent Design Trends and Security Notions. Cryptology ePrint Archive, Report 2011/565. 2011. Available online: <https://eprint.iacr.org/2011/565> (accessed on 1 June 2021).
10. Yoo, H.; Ko, N. Blockchain based Data Marketplace System. In Proceedings of the 2020 International Conference on Information and Communication Technology Convergence (ICTC), Jeju Island, Korea, 21–23 October 2020; pp. 1255–1257. [\[CrossRef\]](#)
11. Mikkelsen, L.; Mortensen, K.; Rasmussen, H.; Schwefel, H.P.; Madsen, T. Realization and Evaluation of Marketplace Functionalities Using Ethereum Blockchain. In Proceedings of the 2018 International Conference on Internet of Things, Embedded Systems and Communications (IINTEC), Hammamet, Tunisia, 20–22 December 2018; pp. 47–52. [\[CrossRef\]](#)
12. Ranganathan, V.P.; Dantu, R.; Paul, A.; Mears, P.; Morozov, K. A Decentralized Marketplace Application on the Ethereum Blockchain. In Proceedings of the 2018 IEEE 4th International Conference on Collaboration and Internet Computing (CIC), Philadelphia, PA, USA, 18–20 October 2018; pp. 90–97. [\[CrossRef\]](#)
13. Braud, A.; Fromentoux, G.; Radier, B.; Le Grand, O. The Road to European Digital Sovereignty with Gaia-X and IDSA. *IEEE Netw.* **2021**, *35*, 4–5. [\[CrossRef\]](#)
14. Özyilmaz, K.R.; Doğan, M.; Yurdakul, A. IDMoB: IoT Data Marketplace on Blockchain. In Proceedings of the 2018 Crypto Valley Conference on Blockchain Technology (CVCBT), Zug, Switzerland, 20–22 June 2018; pp. 11–19. [\[CrossRef\]](#)
15. Nguyen, D.D.; Ali, M.I. Enabling On-Demand Decentralized IoT Collectability Marketplace using Blockchain and Crowdsensing. In Proceedings of the 2019 Global IoT Summit (GIoTS), Aarhus, Denmark, 17–21 June 2019; pp. 1–6. [\[CrossRef\]](#)
16. Tzianos, P.; Pipelidis, G.; Tsiamitros, N. Hermes: An Open and Transparent Marketplace for IoT Sensor Data over Distributed Ledgers. In Proceedings of the 2019 IEEE International Conference on Blockchain and Cryptocurrency (ICBC), Seoul, Korea, 14–17 May 2019; pp. 167–170. [\[CrossRef\]](#)
17. Arya, V.; Sen, S.; Kodeswaran, P. Blockchain Enabled Trustless API Marketplace. In Proceedings of the 2020 International Conference on COMMunication Systems NETWORKS (COMSNETS), Bengaluru, India, 7–11 January 2020; pp. 731–735. [\[CrossRef\]](#)
18. Musso, S.; Perboli, G.; Rosano, M.; Manfredi, A. A Decentralized Marketplace for M2M Economy for Smart Cities. In Proceedings of the 2019 IEEE 28th International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE), Napoli, Italy, 12–14 June 2019; pp. 27–30. [\[CrossRef\]](#)
19. Ramachandran, G.S.; Radhakrishnan, R.; Krishnamachari, B. Towards a Decentralized Data Marketplace for Smart Cities. In Proceedings of the 2018 IEEE International Smart Cities Conference (ISC2), Kansas City, MI, USA, 16–19 September 2018; pp. 1–8. [\[CrossRef\]](#)
20. Jeong, B.G.; Youn, T.Y.; Jho, N.S.; Shin, S.U. Blockchain-Based Data Sharing and Trading Model for the Connected Car. *Sensors* **2020**, *20*, 3141. [\[CrossRef\]](#) [\[PubMed\]](#)
21. Li, Y.N.; Feng, X.; Xie, J.; Feng, H.; Guan, Z.; Wu, Q. A decentralized and secure blockchain platform for open fair data trading. *Concurr. Comput. Pract. Exp.* **2020**, *32*, e5578. [\[CrossRef\]](#)
22. Ma, S.; Mu, Y.; Susilo, W. A Generic Scheme of plaintext-checkable database encryption. *Inf. Sci.* **2018**, *429*, 88–101. [\[CrossRef\]](#)
23. Dai, W.; Dai, C.; Choo, K.K.R.; Cui, C.; Zou, D.; Jin, H. SDTE: A Secure Blockchain-Based Data Trading Ecosystem. *IEEE Trans. Inf. Forensics Secur.* **2020**, *15*, 725–737. [\[CrossRef\]](#)
24. Standardization Initiative, P. JSON for Linking Data. 2021. Available online: <https://json-ld.org/> (accessed on 1 June 2021).
25. Muñoz, J.; Forne, J.; Esparza, O. Certificate revocation system implementation based on the Merkle hash tree. *Int. J. Inf. Secur.* **2004**, *2*, 110–124. [\[CrossRef\]](#)
26. Mansour, T.; Schork, M. *Commutation Relations, Normal Ordering, and Stirling Numbers*; Discrete Mathematics and Its Applications; Chapman and Hall/CRC: Boca Raton, FL, USA, 2015.

## Short Biography of Authors



**Rafael Genés-Durán** is currently a PhD candidate of the Information Security Group (ISG) doing research in distributed ledger technologies and zero-knowledge proofs at Universitat Politècnica de Catalunya. He holds a B.S. degree in Telecommunications Engineering (2017) and a Master in Informatics Engineering (2019). Contact him at [rafael.genes@upc.edu](mailto:rafael.genes@upc.edu).



**Juan Hernández-Serrano** is an associate professor of the Department of Network Engineering of the Universitat Politècnica de Catalunya in Spain, and a researcher of the Information Security Group (ISG). He holds an M.S. in Network Engineering (2002) and a Ph.D. in the field of information security (2008). His research interests has been focused on different aspects of networks security and privacy, including IoT, distributed ledgers, cognitive radio networks, M2M, smart grids, eVoting and digital forensics. Contact him at [j.hernandez@upc.edu](mailto:j.hernandez@upc.edu).



**Oscar Esparza** is working as associate professor of the Department of Network Engineering of the Universitat Politècnica de Catalunya. He holds an M.S. in Telecommunications Engineering (1999) and a PhD in Security Engineering (2004). His expertise areas are related to network security and applied cryptography. Since 2017 he leads the Information Security Group (ISG) of the UPC. Contact him at [oscar.esparza@upc.edu](mailto:oscar.esparza@upc.edu).



**Marta Bellés-Muñoz** received her B.S. degree in Mathematics at Universitat Autònoma de Barcelona and continued her Master studies at Aarhus Universitet, where she focused on the study of elliptic curves and isogeny-based cryptography. She is currently a PhD student doing research on security and efficiency of arithmetic circuits for zero-knowledge proofs at Universitat Pompeu Fabra in collaboration with Dusk Network. Contact her at [marta.belles@upf.edu](mailto:marta.belles@upf.edu).



**Jose L. Muñoz-Tapia** is a researcher of the Information Security Group (ISG) and an associate professor of the Department of Network Engineering of the Universitat Politècnica de Catalunya. He holds an M.S. in Telecommunications Engineering (1999) and a PhD in Security Engineering (2003). He has worked in applied cryptography, network security and game theory models applied to networks and simulators. His research focus has now tuned to distributed ledger technologies, and he is the director of the Master program in Blockchain technologies at UPC School. Contact him at [jose.luis.munoz@upc.edu](mailto:jose.luis.munoz@upc.edu).