

# Payment with Dispute Resolution: A Protocol For Reimbursing Frauds Victims

**Abstract**—An “Authorised Push Payment” (APP) fraud refers to a case where fraudsters deceive a victim to make payments to bank accounts controlled by them. The total amount of money stolen via APP frauds is swiftly growing. Although regulators have provided guidelines to improve victims’ protection, the guidelines are vague, the implementation is lacking in transparency, and the victims are not receiving sufficient protection. To *facilitate victims’ reimbursement*, in this work, we propose a protocol called “Payment with Dispute Resolution” (PwDR) and formally define it. The protocol lets an honest victim prove its innocence to a third-party dispute resolver while preserving the protocol participants’ *privacy*. It makes black-box use of a standard online banking system. We implement its most computationally-intensive subroutine and analyse its runtime. We also evaluate its asymptotic cost. Our evaluation indicates that the protocol is efficient. It imposes only  $O(1)$  overheads to the customer and bank. Moreover, it takes a dispute resolver only 0.09 milliseconds to settle a dispute between the two parties.

## 1. Introduction

An “Authorised Push Payment” (APP) fraud is a type of cyber-crime where a fraudster tricks a victim into making an authorised online payment into an account controlled by the fraudster (Appendix A presents further definitions of APP fraud). The APP fraud has various variants, such as romance, investment, or invoice fraud [48]. The total amount of money lost to APP fraud is substantial. According to statistics collected from the UK banking industry by “UK Finance”, in the first half of 2021, a total of £355 million was lost to APP frauds. Losses have increased by 71% compared to those reported in the same period in 2020 and now makes up almost half of banking fraud losses in the UK [47]. APP fraud is a *global* phenomenon. According to the FBI, victims of APP fraud reported to it at least a total of \$419 million losses, in 2020 [20]. Recently, Interpol warned its member countries about a variant of APP fraud called investment fraud via dating software [46]. According to Europol’s notice, at least five variants of APP fraud are among the seven most common types of online financial fraud [43].

Although the amount of money lost via APP frauds and the number of cases have been significantly increasing, the victims are not receiving enough protection. In the first half of 2021, only 42% of the stolen funds returned to victims of APP frauds in the UK [47]. Despite the UK’s financial regulators (unlike US and EU) having provided specific reimbursement regulations to financial institutes to improve APP fraud victims’ protection, these regulations are ambiguous and open to interpretation. Also, there

exists no transparent and uniform mechanism via which honest victims can *prove* their innocence. Currently, each bank uses its own ad-hoc (manual) dispute resolution process which is not transparent to customers, regulators, or consumer protection organisations. It is not uniform among all banks and even among those organisations that settle disputes between banks and customers. To date, the APP fraud problem has been overlooked by the information security and cryptography research communities.

In this work, to facilitate the compensation of APP frauds victims for their loss, we propose a protocol called “Payment with Dispute Resolution” (PwDR), present its formal definition, and prove the protocol’s security. The PwDR lets an honest victim (of an APP fraud) independently prove its innocence to a (potentially semi-honest) third-party dispute resolver, in order to be reimbursed. We identify three crucial properties that such a scheme should possess; namely, (a) security against a malicious victim: a malicious victim who is not qualified for the reimbursement should not be reimbursed, (b) security against a malicious bank: a malicious bank cannot disqualify an honest victim from being reimbursed, and (c) privacy: the customer’s and bank’s messages remain confidential from non-participants of the scheme, and a party which resolves dispute learns as little information as possible. The PwDR makes black-box use of a standard online banking system, hence it can rely on and extend the security of the existing banking system. It automates the implementation of reimbursement regulations where possible and distributes the role of making more subjective decisions among multiple auditors.

The PwDR offers *transparency* by (i) accurately formalising reimbursements’ conditions: it presents an accurate publicly available formalisation capturing the circumstances under which a customer is reimbursed, (ii) offering traceability: it lets parties’ performance be tracked, and (iii) providing an evidence-based final decision: it requires the reasons leading to the final decision to be accessible and consistent with the reimbursements’ conditions and parties’ actions. It also offers *accountability*, as it is equipped with auditing mechanisms that help identify the party liable for an APP fraud loss. The auditing mechanisms themselves are accompanied by our *novel lightweight privacy-preserving* threshold voting protocols, which let auditors vote privately without having to worry about being retaliated against, for their votes. Our voting protocols can be of independent interest. We analyse the PwDR’s cost via both asymptotic and runtime evaluation. The evaluation indicates that the protocol is indeed efficient. The customer’s and bank’s complexity is constant,  $O(1)$ . It only takes 0.09 milliseconds for a dispute resolver to settle a dispute between the two parties. We make the

implementation source code publicly available. We hope that our result lays the foundation for future solutions that will protect victims of this concerning type of fraud.

In summary, our contributions are three-fold, we (1) propose an efficient protocol called Payment with Dispute Resolution (PwDR), (2) formally define and prove the PwDR, and (3) analyse its asymptotic and concrete costs.

## 2. Background

Losses resulting from APP fraud fall outside legislation that protects customers from “unauthorised” payments (e.g., the Payment Services Directive 2 in the EU and the Electronic Fund Transfer Act in the US). Liability for APP frauds has largely remained with victims who authorise payments. In the UK, there have been efforts to protect the victims. In 2016, the UK’s consumer protection organisation, called “Which?”, submitted a super-complaint to the “Financial Conduct Authority” (FCA) and raised its concerns that despite the APP fraud rate is growing, victims do not have enough protection [22]. Since then, the FCA has been collaborating with financial institutions to develop several initiatives that could improve the response when frauds occur. As a result, the “Contingent Reimbursement Model” (CRM) code [31] was proposed. It lays out a set of requirements and explains under which circumstances customers should be reimbursed by their financial institutions when they fall victim to an APP fraud, e.g., when customers were not grossly negligent, did not ignore effective warnings, or were among a vulnerable group. So far, at least nine firms, comprising nineteen brands (e.g., Barclays, HSBC, Lloyds) signed up to the code.

Although the CRM code is a vital guideline for protecting the victims of fraud, some aspects of it are vague and open to interpretation. For instance, in 2020, the “Financial Ombudsman Service” (which settles complaints between consumers and businesses) had an “overall impression” that firms are applying the CRM code inconsistently and in some cases incorrectly, which resulted in failing to reimburse victims in cases anticipated by this code [45].

Unlike the UK which has already recognised APP frauds and developed regulations for that, the US’s financial industry has not distinguished between APP fraud and other types of fraud when labelling their overall fraud losses [1]. Thus, there exist no publicly available accurate reports and regulations concerning APP frauds provided by the financial industry in the US yet. To address the issue, in 2020, the Federal Reserve introduced the “Fraud-Classifier Model” which classifies payment frauds more granularly [44]. In the same line of effort, very recently (on April 25, 2022) a discussion draft has been proposed in the United States House of Representatives to amend the “Electronic Fund Transfer Act” to treat fraudulently induced electronic fund transfers (a.k.a. APP fraud) in the same manner as unauthorized electronic fund transfer, which will ultimately lead to higher customers protection [34]. It is yet to be seen how regulations related to APP fraud will be developed and how the payment industry will implement these regulations, in the US. Similarly, in the EU, there exists no specific regulation developed to explicitly capture APP frauds and protect the victims [27], [32].

## 3. Related Work

### 3.1. Authorised Push Payment Fraud

Anderson *et al.* [3] provide an overview of APP frauds and highlight that although the (CRM) code would urge banks to accept more liability for APP frauds, it remains to be seen how this will evolve as fraudsters will continuously try to figure out how bank systems can facilitate misdirection attacks. Taylor *et al.* [42] analyse the CRM code from legal and practical perspectives. They state that this code’s proper implementation would make considerable advances to protect victims of APP frauds. Nevertheless, they also argue that the code is still ambiguous. Kjørven [27] investigates whether banks or customers should be liable for customers’ financial loss to online frauds including APP ones, under Scandinavian and European law. The author states that consumers are often left to deal with the losses caused by APP frauds. She concludes that this should change and a larger portion of the losses should be allocated to banks.

### 3.2. Dispute Resolution

In payment platforms, dispute resolution mechanisms can be broadly split into two classes, (a) *centralised* and (b) *decentralised*. In the former class, at any point in time, a *single party* tries to settle a dispute. In particular, if a customer disputes having made or authorised a transaction, then the related bank tries directly settle the dispute with the customer. However, banks’ terms and conditions (T&C) can complicate the dispute resolution process. If they do not reach an agreement, the customer can take its case to a third party (e.g., Financial Ombudsman Service or court) to settle the dispute. In 2000, Bohm *et al.* [10] analysed different terms of banks in the UK. They argued that the approach taken by banks is unfair to their customers in some cases. Later, Anderson [4] points out that the move to online banking led many financial institutions to impose T&C that ultimately would shift the burden of proof in dispute to the customer. Becker *et al.* [8] investigate to what extent bank customers know the T&C they signed up for. Their study suggests that only 35% of customers fully understand T&C and 28% of customers find important parts of T&C are unclear.

Now we turn our attention to the latter class, i.e., decentralised dispute resolution. After the invention of the blockchain technology and especially its vital side-product, smart contract, researchers considered the possibility of resolving disputes in a decentralised manner by relying on smart contracts. Such a possibility has been discussed and studied by the law research community, e.g., in [13], [35], [36]. Moreover, various ad-hoc blockchain-based cryptographic protocols have been proposed to resolve disputes in different contexts and settings. We briefly explain a few of them. Dziembowski *et al.* [17] propose FairSwap, an efficient protocol that allows a seller and buyer to fairly exchange digital items and coins. It is mainly based on a Merkle tree and Ethereum smart contracts which can efficiently resolve a dispute between the seller and buyer when the two parties disagree. Recently, researchers in [18] propose OPTISWAP that improves FairSwap’s performance. Similar to FairSwap,

OPTISWAP uses a Merkle tree and smart contract, but it relies on an *interactive* dispute resolution protocol.

Recently, Abadi *et al.* [2] propose a privacy-preserving protocol that allows a fair exchange of digital coins and a certain *digital service*, called “proofs of data retrievability”. To efficiently settle disputes between a seller and a buyer, the protocol uses blockchain, Merkle tree, symmetric key encryption, and a third-party arbiter. In general, in fair exchange protocols, parties want to “exchange” digital items, e.g., coins and files/services. However, in various variants of APP fraud (e.g., CEO, romance, or invoice) victims do not necessarily expect to receive an item in exchange for the money they transfer. Thus, fair exchange protocols cannot fully solve the problem that we want to solve. However, there are very specific cases (e.g., Purchase APP fraud) where fair exchange protocols can be used to deal with APP fraud.

In the context of verifiable (cloud) computation, Dong *et al.* [16] use a combination of smart contracts, game theory (incentivization), and a third-party arbiter to design an efficient protocol that lets a client outsource its expensive computation to the cloud servers such that it can efficiently check the result’s correctness. In the case of dispute, the protocol lets the parties invoke the arbiter which efficiently settles the dispute with the assistance of the smart contract. Green *et al.* [24] propose a variant of payment channel [39] (which improves cryptocurrencies’ scalability) while preserving the users’ anonymity. In this scheme, in the case of a dispute between two parties, they can send a set of proofs to a smart contract that settles the disputes between the two.

Thus, although many dispute resolution solutions have been proposed, to date, no (centralised or decentralised) solution exists to resolve disputes in the context of APP frauds. Our PwDR is the first protocol that fills in the gap.

## 4. Preliminaries

### 4.1. Notations

We summarise our notations in Table 1.

### 4.2. Informal Threat Model and Assumptions

A payment with a dispute resolution scheme involves six types of parties. Below, we informally explain each type of party’s role. We will provide a formal definition of the scheme in Section 6.

- Customer ( $\mathcal{C}$ ): it is a regular customer of a bank. We call a customer a victim after it falls victim to an APP fraud. We assume a victim is corrupted by a non-colluding active (or malicious) adversary.
- Bank ( $\mathcal{B}$ ): it is a regular bank providing online banking. We assume it is corrupted by a non-colluding active adversary. We assume any change to the online banking system’s source code is transparent and can be detected.
- Smart contract ( $\mathcal{S}$ ): it is a standard smart contract of a public blockchain (e.g., Ethereum). It mainly acts as a tamper-proof public bulletin board to store different parties’ messages.
- Certificate generator ( $\mathcal{G}$ ): it is a trusted third party (e.g., registry office) which provides signed digital certificates (e.g., certificate of disability, divorce) to customers.

TABLE 1: Notation Table.

Symbol	Description
$\text{Enc}(\cdot)$	Encryption algorithm of symmetric key encryption
$\text{Dec}(\cdot)$	Decryption algorithm of symmetric key encryption
$\text{Enc}(\cdot)$	Encryption algorithm of asymmetric key encryption
$\text{Dec}(\cdot)$	Decryption algorithm of asymmetric key encryption
$\text{keyGen}(\cdot)$	Key generator algorithm of asymmetric key encryption
$\text{Sig.keyGen}(\cdot)$	Key generator algorithm of digital signature scheme
$\text{verStat}(\cdot)$	Algorithm to determine $\mathcal{B}$ ’s message status
$\text{checkWarning}(\cdot)$	Algorithm to check a warning’s effectiveness
$\text{pay}(\cdot)$	$\mathcal{B}$ ’s internal algorithm to transfers money
$\text{Com}(\cdot)$ and $\text{Ver}(\cdot)$	Commitment’s commit and verify resp.
$\mathcal{C}$	Customer
$\mathcal{B}$	Bank
$\mathcal{D}_1, \dots, \mathcal{D}_n$	Auditors
$\mathcal{DR}$	Dispute resolver
$\mathcal{S}$	Smart contract
$\mathcal{G}$	Certificate generator
PPT	Probabilistic polynomial time
PVE	Private verdict encoding protocol
GPVE	Generic private verdict encoding protocol
GFVD	Generic final verdict decoding protocol
$z_1$	$\mathcal{C}$ ’s complaint about $\mathcal{B}$ ’s message status
$z_2$	$\mathcal{C}$ ’s complaint about a warning’s effectiveness
$z_3$	$\mathcal{C}$ ’s complaint about payment message inconsistency
$\text{aux}, \text{aux}'$	Auxiliary information
$\Delta$	Time parameter
$\hat{l}$	$\mathcal{C}$ ’s encoded payees list
$f$	New payee’s detail
$\text{in}_f$	Payment detail
$\hat{a}$	Encoded $a$
$\hat{m}_1^{(\mathcal{C})}$	$\mathcal{C}$ ’s encoded update request
$\hat{m}_2^{(\mathcal{C})}$	$\mathcal{C}$ ’s encoded payment request
$\hat{m}_1^{(\mathcal{B})}$	$\mathcal{B}$ ’s encoded warning message
$\hat{m}_2^{(\mathcal{B})}$	$\mathcal{B}$ ’s encoded payment message
$sk$ and $pk$	Secret and public keys
$sk_{\mathcal{D}}$	Auditors’ secret key
$T, T_1, T_2$	Tokens, where $T := (T_1, T_2)$
$w_1$	Output of $\text{verStat}(\cdot)$
$(w_2, w_3)$	Output of $\text{checkWarning}(\cdot)$
$w_4$	If set 1, the payment was indeed made
$\hat{w}_j$	Output of PVE( $\cdot$ )
$v$	Output of FVD( $\cdot$ )
$e$	Threshold
$w_{i,j}$	Auditor’s plain verdict
$o$	Counter
$r_j$	Pseudorandom value
$\lambda$	Security parameter
$\pi$	Private statement
$\text{in}_p$	The input of $\text{pay}(\cdot)$
$\phi$	Null
$n$	Total number of auditors
PwDR	payment with dispute resolution protocol
$\mu(\cdot)$	Negligible function
PRF( $\cdot$ )	Pseudorandom function
$\bar{k}_0$	A key for PRF
$\bar{k}_1, \bar{k}_2$	Keys for the symmetric key encryption

- A committee of auditors ( $\mathcal{D}_1, \dots, \mathcal{D}_n$ ): it consists of trusted third-party authorities or regulators (e.g., FCA, financial ombudsman service). They compile complaints and provide their verdicts. We assume they interacted with each other once, to agree on a secret key,  $\bar{k}_0$ , and a pair of keys  $(pk_{\mathcal{D}}, sk_{\mathcal{D}})$  of an asymmetric key encryption.
- Dispute resolver ( $\mathcal{DR}$ ): it is an aggregator of auditors’ votes (e.g., public court). Given a collection of votes, it extracts and announces the final verdict. We assume it is corrupted by a non-colluding passive adversary. We assume  $\mathcal{C}$  and  $\mathcal{B}$  use a secure channel when they send a message directly to  $\mathcal{DR}$ .

### 4.3. Digital Signature

A digital signature is a scheme for verifying the authenticity of digital messages and is formally defined in [26] as below.

**Definition 1.** A signature scheme involves three algorithms,  $(\text{Sig.keyGen}, \text{Sig.sign}, \text{Sig.ver})$ , that are defined as follows. (1)  $\text{Sig.keyGen}(1^\lambda) \rightarrow (sk, pk)$  is a probabilistic algorithm run by a signer. It takes as input a security parameter. It outputs a key pair:  $(sk, pk)$ , consisting of secret key  $sk$ , and public key  $pk$ . (2)  $\text{Sig.sign}(sk, pk, u) \rightarrow sig$  is an algorithm run by the signer. It takes as input key pair:  $(sk, pk)$  and a message:  $u$ . It outputs a signature:  $sig$ . (3)  $\text{Sig.ver}(pk, u, sig) \rightarrow h \in \{0, 1\}$  is an algorithm run by a verifier. It takes as input public key:  $pk$ , message:  $u$ , and signature:  $sig$ . It checks the signature's validity. If the verification passes, then it outputs 1; otherwise, it outputs 0.

A digital signature scheme must meet two properties: (1) *Correctness*: for every input  $u$  it holds that:  $Pr[\text{Sig.ver}(pk, u, \text{Sig.sign}(sk, pk, u)) = 1 : \text{Sig.keyGen}(1^\lambda) \rightarrow (sk, pk)] = 1$ . And (2) *Existential unforgeability under chosen message attacks*: a probabilistic polynomial time (PPT) adversary that obtains  $pk$  and has access to a signing oracle for messages of its choice, cannot create a valid pair  $(u^*, sig^*)$  for a new message  $u^*$ , except with a small probability,  $\sigma$ . Formally:  $Pr[u^* \notin Q \wedge \text{Sig.ver}(pk, u^*, sig^*) = 1 : \text{Sig.keyGen}(1^\lambda) \rightarrow (sk, pk), \mathcal{A}^{\text{Sig.sign}(sk, \cdot)}(pk) \rightarrow (u^*, sig^*)] \leq \mu(\lambda)$ , where  $Q$  is the set of queries that  $\mathcal{A}$  sent to the oracle.

#### 4.4. Smart Contract

Cryptocurrencies, such as Bitcoin [33] and Ethereum [49], beyond offering a decentralised currency, support computations on transactions. In this setting, a certain computation logic is encoded in a computer program, called a “smart contract”. Although Bitcoin, the first decentralised cryptocurrency, supports smart contracts, the functionality of Bitcoin’s smart contracts is limited. To address this limitation, Ethereum, as a generic smart contract platform, was designed. Thus far, Ethereum has been the most predominant cryptocurrency framework that lets users define arbitrary smart contracts. To prevent a denial-of-service attack, Ethereum requires a transaction creator to pay a fee, called “gas”.

#### 4.5. Commitment Scheme

A commitment scheme involves two parties, *sender* and *receiver*, and includes two phases: *commit* and *open*. In the commit phase, the sender commits to a message  $x$  as  $\text{Com}(x, r) = \text{Com}_x$ , that involves a secret value,  $r$ . In the open phase, the sender sends the opening  $\tilde{x} := (x, r)$  to the receiver which verifies its correctness:  $\text{Ver}(\text{Com}_x, \tilde{x}) \stackrel{?}{=} 1$  and accepts if the output is 1. A commitment scheme satisfies two properties, (a) *hiding*: it is infeasible for an adversary to learn any information about the message, and (b) *binding*: it is infeasible for an adversary to open a commitment to different values than the one used in the commit phase. We provide more detail about the commitment scheme in Appendix B.

#### 4.6. Statement Agreement Protocol

The “Statement Agreement Protocol” (SAP) proposed in [2] lets two mutually distrusted parties, e.g.,  $\mathcal{B}$  and

$\mathcal{C}$ , efficiently agree on a private statement,  $\pi$ . The SAP satisfies four properties: (1) neither party can convince a third-party verifier that it has agreed with its counterparty on a different statement than the one both parties previously agreed on, (2) after they agree on a statement, an honest party can (almost) always prove to the verifier that it has the agreement, (3) the privacy of the statement is preserved (from the public), and (4) after both parties reach an agreement, neither can deny it. It assumes that each party has a blockchain public address,  $adr_{\mathcal{R}}$  (where  $\mathcal{R} \in \{\mathcal{B}, \mathcal{C}\}$ ). Below, we restate the SAP.

1) **Initiate.**  $\text{SAP.init}(1^\lambda, adr_{\mathcal{B}}, adr_{\mathcal{C}}, \pi)$ .

The following steps are taken by  $\mathcal{B}$ .

- Deploys a smart contract that states both parties’ addresses,  $adr_{\mathcal{B}}$  and  $adr_{\mathcal{C}}$ . Let  $adr_{\text{SAP}}$  be the deployed contract’s address.
- Picks a random value  $r$ , and commits to  $\pi$  as  $\text{Com}(\pi, r) = g_{\mathcal{B}}$ . It sends  $adr_{\text{SAP}}$  and  $\tilde{\pi} := (\pi, r)$  to  $\mathcal{C}$ , and  $g_{\mathcal{B}}$  to the contract.

2) **Agreement.**  $\text{SAP.agree}(\pi, r, g_{\mathcal{B}}, adr_{\mathcal{B}}, adr_{\text{SAP}})$ .

The following steps are taken by  $\mathcal{C}$ .

- Checks if  $g_{\mathcal{B}}$  was sent from  $adr_{\mathcal{B}}$ , and checks locally  $\text{Ver}(g_{\mathcal{B}}, \tilde{\pi}) = 1$ .
- If the checks pass, it sets  $b = 1$ , computes locally  $\text{Com}(\pi, r) = g_{\mathcal{C}}$ , and sends  $g_{\mathcal{C}}$  to the contract. Else, it sets  $b = 0$  and  $g_{\mathcal{C}} = \perp$ .

3) **Prove.** For either  $\mathcal{B}$  or  $\mathcal{C}$  to prove, it sends  $\tilde{\pi} := (\pi, r)$  to the smart contract.

4) **Verify.**  $\text{SAP.verify}(\tilde{\pi}, g_{\mathcal{B}}, g_{\mathcal{C}}, adr_{\mathcal{B}}, adr_{\mathcal{C}})$ .

The following steps are taken by the smart contract.

- Ensures  $g_{\mathcal{B}}$  and  $g_{\mathcal{C}}$  were sent from  $adr_{\mathcal{B}}$  and  $adr_{\mathcal{C}}$  respectively. It also ensures  $\text{Ver}(g_{\mathcal{B}}, \tilde{\pi}) = \text{Ver}(g_{\mathcal{C}}, \tilde{\pi}) = 1$ .
- Outputs  $s = 1$ , if the checks in steps 4a pass. It outputs  $s = 0$ , otherwise.

#### 4.7. Pseudorandom Function

Informally, a pseudorandom function (as is defined in [26]) is a deterministic function that takes a key of length  $\Lambda$  and an input; it outputs a value indistinguishable from that of a truly random function. In this paper, we use the pseudorandom function:  $\text{PRF} : \{0, 1\}^\Lambda \times \{0, 1\}^* \rightarrow \mathbb{F}_p$ , where  $p$  is a large prime number,  $|p| = \lambda$ , and  $(\Lambda, \lambda)$  are the security parameters. In practice, a pseudorandom function can be obtained from an efficient block cipher, e.g., AES, [26].

#### 4.8. Bloom Filter

A Bloom filter [9] is a compact data structure that lets us efficiently check an element membership. It is an array of  $m$  bits (initially all set to zero), that represents  $n$  elements. It is accompanied by  $k$  independent hash functions. To insert an element, all the hash values of the element are computed and their corresponding bits in the filter are set to 1. To check an element membership, all its hash values are re-computed and checked whether all are set to 1 in the filter. If all the corresponding bits are 1, then the element is probably in the filter; otherwise, it is not. In this work, we require that a Bloom filter uses *cryptographic* hash functions. In Appendix C, we explain how the Bloom filter’s parameters can be set.

## 5. Challenges to Overcome

Our starting point in defining and designing a payment with dispute resolution scheme is the CRM code, as this code (although vaguely) sets out the primary requirements a victim must meet to be reimbursed. To design such a scheme, we need to address several challenges. The rest of this section outlines these challenges.

### 5.1. Challenge 1: Lack of Transparent Logs

In the current online banking system, during a payment journey, the messages exchanged between customer and bank are usually logged by the bank and are not accessible to the customer without the bank's collaboration. Even if the bank provides access to the transaction logs, there is no guarantee that the logs have remained intact.

Due to the lack of a transparent logging mechanism, a customer or bank can wrongly claim that (a) it has sent a certain message or warning to its counter-party or (b) it has never received a certain message. Thus, it would be hard for an honest party to prove its innocence. To address this challenge, our scheme will use a smart contract to which each party sends its messages.

### 5.2. Challenge 2: Lack of Effective Warning's Accurate Definition in Banking

One of the determining factors in the process of allocating liability to an APP fraud victim is following "warning(s)", according to the CRM code. However, there exists no publicly available study on the effectiveness of banks' warnings. So, we cannot hold a customer accountable for becoming a fraud victim, even if the related warnings are ignored. Also, currently, banks assess whether their own warnings are effective. But, in a fair process, such an assessment is conducted by a neutral third party.

To address these challenges, we let a warning's effectiveness be determined on a case-by-case basis after an APP fraud occurs. The protocol lets a victim challenge a certain warning whose effectiveness will be assessed by a *committee*, i.e., a set of auditors. In this setting, each auditor provides its (encoded) verdict to the smart contract, from which a dispute resolver retrieves all verdicts to learn the final one. The scheme ensures that the final verdict is in the customer's favour if at least a threshold of the auditors voted so. Thus, unlike the traditional setting where a central party determines a warning's effectiveness, which is error-prone, we let a collection of auditors determine it.

### 5.3. Challenge 3: Linking Off-chain Payments with a Smart Contract

Recall that an APP fraud occurs when a payment is made. In the case where a bank sends (to the smart contract) a confirmation of payment message, it is not possible to automatically validate such a claim, as the money transfer occurs outside of the blockchain network. To address this challenge, our scheme lets a customer raise a dispute and report it to the smart contract when it detects an inconsistency. In this case, the above auditors investigate and provide their verdicts to the smart contract.

Then, dispute resolver  $\mathcal{DR}$  extracts them and announces the final verdict.

### 5.4. Challenge 4: Preserving Privacy

Although the use of a public logging mechanism is vital in resolving disputes transparently, if it does not use a privacy-preserving mechanism, parties' privacy would be violated. To protect the privacy of the bank's and customers' messages from the public, our scheme lets them provably agree on encoding-decoding tokens with which they can encode their messages.

Later, either party can provide the token to a third party (e.g.,  $\mathcal{D}_i$ ) which checks the token's correctness, and decodes the messages. To protect the privacy of the committee members' verdicts from  $\mathcal{DR}$ , the scheme ensures that  $\mathcal{DR}$  learns only the final verdict without being able to link a verdict to a specific auditor or even learn the number of yes/1 and no/0 votes. To this end, we develop and use novel threshold voting protocols.

## 6. Definition of Payment with Dispute Resolution Scheme

In this section, we outline a formal definition of the payment with dispute resolution (pwwr) notion. We refer readers to Appendix D for the full version of the definition.

**Definition 2.** A pwwr involves six types of entities; namely, bank  $\mathcal{B}$ , customer  $\mathcal{C}$ , smart contract  $\mathcal{S}$ , certificate generator  $\mathcal{G}$ , set of auditors  $\mathcal{D} : \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$ , and dispute resolver  $\mathcal{DR}$ . It also includes the following algorithms.

- $\text{keyGen}(1^\lambda) \rightarrow (sk, pk)$ . It generates and outputs a pair of secret keys  $sk := (sk_{\mathcal{G}}, sk_{\mathcal{D}})$  and public keys  $pk := (pk_{\mathcal{G}}, pk_{\mathcal{D}})$ , where  $sk_{\mathcal{D}}$  may include multiple secret keys.
- $\text{bankInit}(1^\lambda) \rightarrow (T, pp, \mathbf{l})$ . It outputs an encoding-decoding token  $T$  (where  $T := (T_1, T_2)$ , each  $T_i$  contains a secret value  $\pi_i$  and its public witness  $g_i$ ), set of public parameters  $pp$  (including a threshold parameter  $e$ ), and empty list  $\mathbf{l}$ .
- $\text{customerInit}(1^\lambda, T, pp) \rightarrow a$ . It is an initiation algorithm that checks the correctness of the elements in  $T$  and  $pp$ . If the checks pass, it outputs 1. Else, it outputs 0.
- $\text{genUpdateRequest}(T, f, \mathbf{l}) \rightarrow \hat{m}_1^{(c)}$ . It is an update request algorithm. It uses the new payee's detail  $f$  and encoding algorithm  $\text{Encode}(T_1, \cdot)$  to generate an encoded update request  $\hat{m}_1^{(c)}$ . It outputs  $\hat{m}_1^{(c)}$ .
- $\text{insertNewPayee}(\hat{m}_1^{(c)}, \mathbf{l}) \rightarrow \hat{\mathbf{l}}$ . It is an algorithm that inserts a new payee's detail into  $\mathbf{l}$  and outputs an updated list  $\hat{\mathbf{l}}$ .
- $\text{genWarning}(T, \hat{\mathbf{l}}, aux) \rightarrow \hat{m}_1^{(B)}$ . It is a warning generating algorithm that outputs an encoded (warning) message  $\hat{m}_1^{(B)}$ , with the help of auxiliary data  $aux$  and  $\text{Encode}(T_1, \cdot)$ , where the plaintext message is either "pass" or "warning" string.
- $\text{genPaymentRequest}(T, in_f, \hat{\mathbf{l}}, \hat{m}_1^{(B)}) \rightarrow \hat{m}_2^{(c)}$ . It is an algorithm that generates an encoded payment request

$\hat{m}_2^{(c)}$ , with the help of new payment's detail  $in_f$  and  $\text{Encode}(T_1, \cdot)$ . It outputs  $\hat{m}_2^{(c)}$ .

- $\text{makePayment}(T, \hat{m}_2^{(c)}) \rightarrow \hat{m}_2^{(B)}$ . It generates and outputs an encoded message  $\hat{m}_2^{(B)}$  for confirmation of payment.
- $\text{genComplaint}(\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f) \rightarrow (\hat{z}, \hat{\pi})$ . It generates complaints with the help of auxiliary data  $aux_f$ . If  $\mathcal{C}$  wants to complain that (i) “pass” message should have been a warning or (ii) no message was provided, it sets  $z_1$  to “challenge message”. If its complaint is about the warning's effectiveness, it sets  $z_2$  to a combination of an evidence  $u \in aux_f$ , the evidence's certificate  $sig \in aux_f$ , the certificate's public parameter, and “challenge warning”, where the certificate is obtained from  $\mathcal{G}$  via a query,  $Q$ . If its complaint is about the payment, it sets  $z_3$  to “challenge payment”. It generates and outputs (i) encoded complaints  $\hat{z}$  using  $\text{Encode}(T_1, \cdot)$ , and (ii) encoded secret parameters  $\hat{\pi}$  using another encoding algorithm  $\text{Encode}(pk_D, \cdot)$ .
- $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ . It compiles  $j$ -th auditor's complaints. It initially sets parameters as  $w_{1,j} = w_{2,j} = w_{3,j} = w_{4,j} = 0$ . If the complaint in  $z_1$  is valid, it sets  $w_{1,j} = 1$ . If the certificate in  $z_2$  is valid, it sets  $w_{3,j} = 1$ . It checks the warning's effectiveness, by running algorithm  $\text{checkWarning}(\cdot)$ . If it is not effective, i.e.,  $\text{checkWarning}(m_1^{(B)}) = 0$ , it sets  $w_{2,j} = 1$ . Also, if the payment was indeed made, it sets  $w_{4,j} = 1$ . It outputs encoded verdicts  $\hat{w}_j = [\hat{w}_{1,j}, \hat{w}_{2,j}, \hat{w}_{3,j}, \hat{w}_{4,j}]$  for  $j$ -th auditor.
- $\text{resDispute}(T_2, \hat{w}, pp) \rightarrow v$ . It aggregates all encoded verdicts  $\hat{w} = [\hat{w}_1, \dots, \hat{w}_n]$  and outputs  $v = [v_1, \dots, v_4]$ , where  $v_i = 1$  if at least  $e$  verdicts  $w_{i,j}$  is 1; otherwise,  $v_i = 0$ . If  $v_4 = 1$  and (i) either  $v_1 = 1$  (ii) or  $v_2 = 1$  and  $v_3 = 1$ , then  $\mathcal{C}$  is reimbursed.

A pwdr has two properties, *correctness* and *security*. Correctness requires that the payment journey is completed (in the absence of a fraudster) without the need for (i) the honest customer to complain and (ii) the honest bank to reimburse. A pwdr is secure if it meets three main properties, (a) security against a malicious victim, (b) security against a malicious bank, and (c) privacy.

Informally, security against a malicious victim states that an APP fraud victim who is not qualified for the reimbursement should not be reimbursed. Specifically, a corrupt victim cannot (a) make at least the threshold of the auditors,  $\mathcal{D}_j$ s, conclude that  $\mathcal{B}$  should have provided a warning, although  $\mathcal{B}$  has done so, or (b) make  $\mathcal{DR}$  conclude that the pass message was incorrectly given or a vital warning message was missing despite only less than the threshold of  $\mathcal{D}_j$ s believing so, or (c) persuade at least the threshold of  $\mathcal{D}_j$ s to conclude that the warning was ineffective although it was effective, or (d) make  $\mathcal{DR}$  believe that the warning message was ineffective although only less than the threshold of  $\mathcal{D}_j$ s believe it, or (e) convince  $\mathcal{D}_j$ s to accept an invalid certificate, or (f) make  $\mathcal{DR}$  believe that at least the threshold of  $\mathcal{D}_j$ s accepted the certificate although they did not. Below, we formally state it.

**Definition 3** (Security against a malicious victim). A pwdr is secure against a malicious victim, if for any security parameter  $\lambda$ , auxiliary data  $aux$ , and PPT adversary  $\mathcal{A}$ , there is a negligible function  $\mu(\cdot)$ , such that for experiment  $\text{Exp}_1^{\mathcal{A}}$ :

$\text{Exp}_1^{\mathcal{A}}(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow$  ( $sk, pk$ )
bankInit( $1^\lambda$ )  $\rightarrow$  ( $T, pp, l$ )
 $\mathcal{A}(1^\lambda, T, pp, l) \rightarrow \hat{m}_1^{(c)}$ 
insertNewPayee( $\hat{m}_1^{(c)}, l$ )  $\rightarrow \hat{l}$ 
genWarning( $T, \hat{l}, aux$ )  $\rightarrow \hat{m}_1^{(B)}$ 
 $\mathcal{A}(T, \hat{l}, \hat{m}_1^{(B)}) \rightarrow \hat{m}_2^{(c)}$ 
makePayment( $T, \hat{m}_2^{(c)}$ )  $\rightarrow \hat{m}_2^{(B)}$ 
 $\mathcal{A}(\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk) \rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n] :$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v = [v_1, \dots, v_4]$ 

```

it holds that the following probability is negligible (i.e.,  $\mu(\lambda)$ ):

$$\Pr \left[ \begin{aligned} & \left( (m_1^{(B)} = \text{warning}) \wedge \left( \sum_{j=1}^n w_{1,j} \geq e \right) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{1,j} < e \right) \wedge (v_1 = 1) \right) \vee \\ & \left( (\text{checkWarning}(m_1^{(B)}) = 1) \wedge \right. \\ & \left. \left( \sum_{j=1}^n w_{2,j} \geq e \right) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{2,j} < e \right) \wedge (v_2 = 1) \right) \vee \\ & \left( u \notin Q \wedge \text{Sig.ver}(pk, u, sig) = 1 \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{3,j} < e \right) \wedge (v_3 = 1) \right) \end{aligned} \right] : \text{Exp}_1^{\mathcal{A}}(\text{input})$$

where  $\hat{w}_j = [\hat{w}_{1,j}, \hat{w}_{2,j}, \hat{w}_{3,j}, \hat{w}_{4,j}]$ ,  $\hat{w} = [\hat{w}_1, \dots, \hat{w}_n]$ ,  $\hat{m} = [\hat{m}_1^{(c)}, \hat{m}_2^{(c)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ ,  $(w_{1,j}, \dots, w_{3,j})$  are the decoding of  $(\hat{w}_{1,j}, \dots, \hat{w}_{3,j}) \in \hat{w}_j \in \hat{w}$ , and  $\text{input} := (1^\lambda, aux)$ .

Security against a malicious bank requires that a malicious bank cannot disqualify an honest victim from being reimbursed. Specifically, a corrupt bank cannot (a) make  $\mathcal{DR}$  conclude that the “pass” message was correctly given or an important warning was not missing although at least the threshold of  $\mathcal{D}_j$ s do not believe so, or (b) convince  $\mathcal{DR}$  that the warning message was effective although at least the threshold of  $\mathcal{D}_j$ s do not believe so, or (c) make  $\mathcal{DR}$  believe that less than the threshold of  $\mathcal{D}_j$ s did not accept the certificate although at least the threshold of them did it, or (d) make  $\mathcal{DR}$  believe that no payment was made, although at least the threshold of  $\mathcal{D}_j$ s believe the opposite.

**Definition 4** (Security against a malicious bank). A pwdr scheme is secure against a malicious bank, if for any  $\lambda$ ,  $aux$ , and PPT adversary  $\mathcal{A}$ , there exists a negligible function  $\mu(\cdot)$ , such that for an experiment  $\text{Exp}_2^{\mathcal{A}}$ :

$\text{Exp}_2^A(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow (sk, pk)$ 
 $\mathcal{A}(1^\lambda) \rightarrow (T, pp, l, f, in_f, aux_f)$ 
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
genUpdateRequest( $T, f, l$ )  $\rightarrow \hat{m}_1^{(C)}$ 
insertNewPayee( $\hat{m}_1^{(C)}, l$ )  $\rightarrow \hat{l}$ 
 $\mathcal{A}(T, \hat{l}, aux) \rightarrow \hat{m}_1^{(B)}$ 
genPaymentRequest( $T, in_f, \hat{l}, \hat{m}_1^{(B)}$ )  $\rightarrow \hat{m}_2^{(C)}$ 
 $\mathcal{A}(T, \hat{m}_2^{(C)}) \rightarrow \hat{m}_2^{(B)}$ 
genComplaint( $\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f$ )  $\rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n] :$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v = [v_1, \dots, v_4]$ 

```

it holds that the following probability is  $\mu(\lambda)$ :

$$\Pr \left[ \begin{array}{l} \left( \left( \sum_{j=1}^n w_{1,j} \geq e \right) \wedge (v_1 = 0) \right) \vee \\ \left( \left( \sum_{j=1}^n w_{2,j} \geq e \right) \wedge (v_2 = 0) \right) \vee \\ \left( \left( \sum_{j=1}^n w_{3,j} \geq e \right) \wedge (v_3 = 0) \right) \vee \\ \left( \left( \sum_{j=1}^n w_{4,j} \geq e \right) \wedge (v_4 = 0) \right) \end{array} \right] : \text{Exp}_2^A(\text{input})$$

Informally, a pwdr scheme is privacy-preserving if it protects the privacy of (1) customers, bank, and auditors' sensitive messages from the scheme's non-participants and (2) each auditor's verdict from  $\mathcal{DR}$ .

**Definition 5** (Privacy). A pwdr scheme preserves privacy if the following two properties are satisfied.

- 1) For any PPT adversary  $\mathcal{A}_1$ , security parameter  $\lambda$ , and auxiliary information  $aux$ , there exists a negligible function  $\mu(\cdot)$ , such that for any experiment  $\text{Exp}_3^{\mathcal{A}_1}$ :

$\text{Exp}_3^{\mathcal{A}_1}(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow (sk, pk)$ 
bankInit( $1^\lambda$ )  $\rightarrow (T, pp, l)$ 
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
 $\mathcal{A}_1(1^\lambda, pk, a, pp, g, l) \rightarrow ((f_0, f_1), (in_{f_0}, in_{f_1}),$ 
  ( $aux_{f_0}, aux_{f_1}$ ))
 $\gamma \xleftarrow{\$} \{0, 1\}$ 
genUpdateRequest( $T, f_\gamma, l$ )  $\rightarrow \hat{m}_1^{(C)}$ 
insertNewPayee( $\hat{m}_1^{(C)}, l$ )  $\rightarrow \hat{l}$ 
genWarning( $T, \hat{l}, aux$ )  $\rightarrow \hat{m}_1^{(B)}$ 
genPaymentRequest( $T, in_{f_\gamma}, \hat{l}, \hat{m}_1^{(B)}$ )  $\rightarrow \hat{m}_2^{(C)}$ 
makePayment( $T, \hat{m}_2^{(C)}$ )  $\rightarrow \hat{m}_2^{(B)}$ 
genComplaint( $\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_{f_\gamma}$ )  $\rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n] :$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v$ 

```

it holds that:

$$\Pr [\mathcal{A}_1(g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w}) \rightarrow \gamma : \text{Exp}_3^{\mathcal{A}_1}(\text{input})] \leq \frac{1}{2} + \mu(\lambda)$$

- 2) For any PPT adversaries  $\mathcal{A}_2$  and  $\mathcal{A}_3$ , security parameter  $\lambda$ , and auxiliary information  $aux$ , there exists a negligible function  $\mu(\cdot)$ , such that for any experiment  $\text{Exp}_4^{\mathcal{A}_2}$ :

$\text{Exp}_4^{\mathcal{A}_2}(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow (sk, pk)$ 
bankInit( $1^\lambda$ )  $\rightarrow (T, pp, l)$ 
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
 $\mathcal{A}_2(1^\lambda, pk, a, pp, l) \rightarrow (f, in_f, aux_f)$ 
genUpdateRequest( $T, f, l$ )  $\rightarrow \hat{m}_1^{(C)}$ 
insertNewPayee( $\hat{m}_1^{(C)}, l$ )  $\rightarrow \hat{l}$ 
 $\mathcal{A}_2(T, \hat{l}, aux) \rightarrow \hat{m}_1^{(B)}$ 
Encode( $T_1, \hat{m}_1^{(B)}$ )  $\rightarrow \hat{m}_1^{(S)}$ 
genPaymentRequest( $T, in_f, \hat{l}, \hat{m}_1^{(S)}$ )  $\rightarrow \hat{m}_2^{(C)}$ 
 $\mathcal{A}_2(T, pk, aux_f, \hat{m}_1^{(S)}, \hat{m}_2^{(C)}) \rightarrow (\hat{m}_2^{(B)}, z, \hat{\pi})$ 
Encode( $T_1, \hat{m}_2^{(B)}$ )  $\rightarrow \hat{m}_2^{(S)}$ 
Encode( $T_1, z$ )  $\rightarrow \hat{z}$ 
Encode( $pk_D, \hat{\pi}$ )  $\rightarrow \hat{\pi}$ 
 $\forall j, j \in [n] :$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v$ 

```

it holds that:

$$\Pr \left[ \mathcal{A}_3(T_2, pk, pp, g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w}) : \text{Exp}_4^{\mathcal{A}_2}(\text{input}) \right] \leq Pr' + \mu(\lambda)$$

Let auditor  $\mathcal{D}_i$  output 0 and 1 with probabilities  $Pr_{i,0}$  and  $Pr_{i,1}$  respectively. Then,  $Pr'$  is defined as  $\text{Max}(Pr_{1,0}, Pr_{1,1}, \dots, Pr_{n,0}, Pr_{n,1})$ .

**Definition 6** (Security). A pwdr scheme is secure if it meets security against a malicious victim, security against a malicious bank, and preserves privacy with respect to definitions 3, 4, and 5 respectively.

We refer readers to Appendix E, for further discussion about the privacy definition.

## 7. Payment with Dispute Resolution Protocol

In this section, we first present an outline of the PwDR protocol (in Section 7.1). Then, we present a few subroutines (in Sections 7.2–7.4) that will be used in this protocol. After that, we describe the PwDR protocol in detail (in Section 7.5).

### 7.1. An Overview of the Protocol

At a high level, the PwDR works as follows. Initially,  $\mathcal{C}$  and  $\mathcal{B}$  agree on a smart contract  $\mathcal{S}$ . They also use the SAP to agree on two private statements including two secret keys that will be used to encrypt outgoing messages. When  $\mathcal{C}$  wants to transfer money to a new payee, it signs into its online banking. It generates an update request (that specifies the new payee's detail), encrypts it, and sends the result to  $\mathcal{S}$ . Then,  $\mathcal{B}$  decrypts and checks the request, e.g., whether it meets its internal policy. Depending on the request,  $\mathcal{B}$  generates a pass or warning message. It encrypts the message and sends the result to  $\mathcal{S}$ . Next,  $\mathcal{C}$  checks  $\mathcal{B}$ 's message and decides whether to make payment. If it decides to do so, it sends an encrypted payment detail to  $\mathcal{S}$ . After that,  $\mathcal{B}$  decrypts the message and locally transfers the amount of money specified in  $\mathcal{C}$ 's message. Once the money is transferred,  $\mathcal{B}$  sends an encrypted "paid" message to  $\mathcal{S}$ .

Once  $\mathcal{C}$  realises that it has fallen victim, it raises a dispute. Specifically, it generates an encrypted complaint

that can challenge the effectiveness of the warning and/or any payment inconsistency. It can include in the complaint an evidence/certificate, e.g., asserting that it falls into the vulnerable customer category as defined in the CRM code.  $\mathcal{C}$  encrypts the complaint and sends to  $\mathcal{S}$  the result and a proof asserting the secret key's correctness. Then, each auditor verifies the proof. If the verification passes, it decrypts and compiles  $\mathcal{C}$ 's complaint to generate a (set of) verdict. Each auditor encodes its verdict and sends the encoded verdict's encryption to  $\mathcal{S}$ . To resolve a dispute between  $\mathcal{C}$  and  $\mathcal{B}$ , either of them invokes  $\mathcal{DR}$ . To do so, it directly sends to  $\mathcal{DR}$  one of the above secret keys and a proof asserting that the key was generated correctly.  $\mathcal{DR}$  verifies the proof. If approved, it locally decrypts the encrypted encoded verdicts (retrieved from  $\mathcal{S}$ ) and finds out the final verdict. If the final verdict indicates the legitimacy of  $\mathcal{C}$ 's complaint, then  $\mathcal{C}$  must be reimbursed. Note, the verdicts are encoded in a way that even after decrypting them,  $\mathcal{DR}$  cannot link a verdict to a committee member or even figure out how many 1 or 0 verdicts were provided (except when all verdicts are 0). However, it can find out whether at least the threshold of the auditors voted in favour of  $\mathcal{C}$ .

## 7.2. A Subroutine for Determining Bank's Message Status

In the payment journey, the customer may receive a "pass" message or even nothing at all, e.g., due to a system failure. In such cases, a victim must be able to complain that if the pass or missing message was a warning, then it would have prevented it from falling victim. To assist the auditors to deal with such complaints deterministically, we propose  $\text{verStat}(\cdot)$  algorithm, which is run locally by each committee member. This algorithm is presented in Figure 1.

**verStat**( $add_S, m^{(B)}, l, \Delta, aux$ )  $\rightarrow w_1$

- **Input.**  $add_S$ : the address of smart contract  $\mathcal{S}$ ,  $m^{(B)}$ :  $\mathcal{B}$ 's warning message,  $l$ : customer's payees' list,  $\Delta$ : a time parameter, and  $aux$ : auxiliary information, e.g., bank's policy.
- **Output.**  $w_1 = 0$ : if the "pass" message had been given correctly or the missing message did not play any role in preventing the fraud;  $w_1 = 1$ : otherwise.

- 1) reads the content of  $\mathcal{S}$ . It checks if  $m^{(B)}$  = "pass" or the encrypted warning message was not sent on time (i.e., never sent or sent after  $t_0 + \Delta$ ). If one of the checks passes, it proceeds; Otherwise, it aborts.
- 2) checks the validity of customer's most recent payees' list  $l$ , with the help of  $aux$ .
  - if  $l$  contains an invalid element, it sets  $w_1 = 1$ .
  - otherwise, it sets  $w_1 = 0$ .
- 3) returns  $w_1$ .

Figure 1: Algorithm to Determine a Bank's Message Status.

## 7.3. A Subroutine for Checking a Warning's Effectiveness

To help the auditors deterministically compile a victim's complaint about a warning's effectiveness, we propose an algorithm, called  $\text{checkWarning}(\cdot)$  which is run

locally by each auditor. It also allows the victims to provide (to the auditors) a certificate/evidence as part of their complaints. This algorithm is presented in Figure 2.

**checkWarning**( $add_S, z, m^{(B)}, aux'$ )  $\rightarrow (w_2, w_3)$

- **Input.**  $add_S$ : the address of smart contract  $\mathcal{S}$ ,  $z$ :  $\mathcal{C}$ 's complaint,  $m^{(B)}$ :  $\mathcal{B}$ 's warning message, and  $aux'$ : auxiliary information, e.g., guideline on warnings' effectiveness.
- **Output.**  $w_2 = 0$ : if the given warning message is effective;  $w_2 = 1$ : if the warning message is ineffective. Also,  $w_3 = 1$ : if the certificate in  $z$  is valid or no certificate is provided;  $w_3 = 0$ : if the certificate is invalid.

- 1) parse  $z = m || sig || pk || \text{"challenge warning"}$ . If  $sig$  is empty, it sets  $w_3 = 0$  and goes to step 2. Otherwise, it:
  - a) verifies the certificate:  $\text{Sig.ver}(pk, m, sig) \rightarrow h$ .
  - b) if the certificate is rejected (i.e.,  $h = 0$ ), it sets  $w_3 = 0$ . It goes to step 4.
  - c) otherwise (i.e.,  $h = 1$ ), it sets  $w_3 = 1$  and moves onto the next step.
- 2) checks if "warning"  $\in m^{(B)}$ . If the check is passed, it proceeds to the next step. Otherwise, it aborts.
- 3) checks the warning's effectiveness, with the assistance of the evidence  $m$  and auxiliary information  $aux'$ .
  - if it is effective, it sets  $w_2 = 0$ . Otherwise, it sets  $w_2 = 1$ .
- 4) returns  $(w_2, w_3)$ .

Figure 2: Algorithm to Check Warning's Effectiveness.

## 7.4. Subroutines for Encoding-Decoding Verdicts

In this section, we present verdict encoding and decoding protocols. They let a third party  $\mathcal{I}$ , e.g.,  $\mathcal{DR}$ , learn if a threshold of the auditors voted 1, while satisfying the following requirements. The protocols (1) generate unlinkable verdicts, (2) do not require auditors to interact with each other for each customer, and (3) are efficient. Since the second and third requirements are self-explanatory, we only explain the first one. Informally, the first property states that the protocols generate encoded verdicts and final verdict in a way that  $\mathcal{I}$ , given these values, cannot (a) link a verdict to an auditor (except when all verdicts are 0), and (b) learn the total number of 1 or 0 verdicts when they provide different verdicts. Shortly, we present two variants of verdict encoding and decoding protocol. The first variant is highly efficient and suitable when the threshold is 1. The second one is generic and works for any threshold (but is less efficient).

**7.4.1. Variant 1: Efficient Verdict Encoding-Decoding Protocol.** This variant has two protocols, Private Verdict Encoding (PVE) and Final Verdict Decoding (FVD). They let  $\mathcal{I}$  learn if at least one auditor voted 1. This variant relies on our observation that if a set of random values and 0s are XORed, then the result reveals nothing, e.g., about the number of non-zero and zero values. In Appendix F, we present the above observation's formal statement and its proof. At a high level, PVE and FVD work as follows. The auditors only once agree on a secret key. This key will let each of them, in PVE, generate a pseudorandom masking value such that if all masking values are XORed,



they would cancel out each other.<sup>1</sup> In PVE, each auditor encodes its verdict by (i) representing it as a parameter which is set to 0 if the verdict is 0, or to a random value if the verdict is 1, and then (ii) “masking” this parameter with the above pseudorandom value. It sends the result to  $\mathcal{I}$ . In FVD,  $\mathcal{I}$  XORs all encoded verdicts. This removes the masks and XORs all verdicts’ representations. If the result is 0, it concludes that all auditors voted 0; so, the final verdict is 0. But, if the result is not 0, it knows that at least one of the auditors voted 1, so the final verdict is 1. Figures 3 and 4 present PVE and FVD respectively.

**7.4.2. Variant 2: Generic Verdict Encoding-Decoding Protocol.** This variant also includes two protocols, Generic Private Verdict Encoding (GPVE) and Generic Final Verdict Decoding (GFVD) which let  $\mathcal{I}$  learn if at least  $e$  auditors voted 1, where  $e$  is an integer in  $[1, n]$ . It uses a novel combination of Bloom filter and combinatorics. It relies on our observation that a Bloom filter encoding a set of random values reveals nothing about the set’s elements. Appendix H presents the above observation’s formal statement and proof. In this variant also, the auditors initially agree on a secret key used to generate a pseudorandom masking value. Each auditor  $\mathcal{D}_j$  represents its verdict by a parameter, such that if its verdict is 0, it sets the parameter to 0; but, if the verdict is 1, it sets the parameter to a fresh *pseudorandom* value  $\alpha_j$ , also derived from the above key. Thus, there would be a set  $A = \{\alpha_1, \dots, \alpha_n\}$  from which  $\mathcal{D}_j$  would pick  $\alpha_j$  to represent its verdict 1.

Each  $\mathcal{D}_j$  masks its verdict representation by its masking value. It sends the result to  $\mathcal{I}$ . Also, (only) auditor  $\mathcal{D}_n$  generates a set  $W$  of all combinations of auditors’ verdict 1’s representations that satisfy the threshold,  $e$ . Specifically, for every integer  $i$  in  $[e, n]$ , it computes the combinations (without repetition) of  $i$  elements from  $A = \{\alpha_1, \dots, \alpha_n\}$ . If multiple elements are taken at a time (i.e.,  $i > 1$ ), they are XORed with each other. Let  $W = \{(\alpha_1 \oplus \dots \oplus \alpha_e), (\alpha_2 \oplus \dots \oplus \alpha_{e+1}), \dots, (\alpha_1 \oplus \dots \oplus \alpha_n)\}$  be the result.  $\mathcal{D}_n$  computes each element of  $W$  regardless of what a specific auditor votes; also, it can generate each  $\alpha_i$  independently (without interacting with other auditors), as it knows the single key (of the pseudorandom function) that was used by other auditors to generate these values. To protect the votes representations’ privacy (from  $\mathcal{I}$ ), it inserts all elements of  $W$  into a Bloom filter. Let BF be the resulting Bloom filter. It sends BF to  $\mathcal{I}$ .

In GFVD, to decode and extract the final verdict,  $\mathcal{I}$  XORs all masked verdict representations which removes the masking values and XORs the representations. Let  $c$  be the result. If  $c = 0$ , then  $\mathcal{I}$  concludes that all auditors voted 0; so, it sets the final verdict to 0. If  $c \neq 0$ , then it checks if  $c \in \text{BF}$ . If it is, then it concludes that at least the threshold of the auditors voted 1, so it sets the final verdict to 1. Otherwise ( $c \notin \text{BF}$ ), it learns that less than the threshold of the auditors voted 1; so, it sets the final verdict to 0. Figures 6 and 7, in Appendix G, present the GPVE and GFVD protocols in detail. Note, the total number of the combinations, i.e., the cardinality of  $W$ , is small when the number of auditors is not very high. In

general, due to the binomial theorem, the cardinality of  $W$  is determined as:  $|W| = \sum_{i=e}^n \frac{n!}{i!(n-i)!}$

For instance, when  $n = 10$  and  $e = 6$ , then  $|W|$  is only 386. Appendix I provides further discussion on the above protocols.

**PVE( $\bar{k}_0, \text{ID}, w_j, o, n, j$ )  $\rightarrow \bar{w}_j$**

- **Input.**  $\bar{k}_0$ : a key of pseudorandom function  $\text{PRF}(\cdot)$ , ID: a unique identifier,  $w_j$ : a verdict,  $o$ : a counter,  $n$ : the total number of auditors, and  $j$ : an auditor’s index.
- **Output.**  $\bar{w}_j$ : an encoded verdict.

Auditor  $\mathcal{D}_j$  takes the following steps.

1) computes a pseudorandom value, as follows.

- if  $j < n$ :  $r_j = \text{PRF}(\bar{k}_0, o || j || \text{ID})$ .

- if  $j = n$ :  $r_j = \bigoplus_{i=1}^{n-1} r_i$ .

2) sets a fresh parameter,  $w'_j$ , as below.

$$w'_j = \begin{cases} 0, & \text{if } w_j = 0 \\ \alpha_j \xleftarrow{\$} \mathbb{F}_p, & \text{if } w_j = 1 \end{cases}$$

3) encodes  $w'_j$  as follows.  $\bar{w}_j = w'_j \oplus r_j$ .

4) outputs  $\bar{w}_j$ .

Figure 3: Private Verdict Encoding (PVE) Protocol. In the figure,  $\mathcal{D}_n$  can generate other auditors’  $r_i$  values, given  $\bar{k}_0$ . Note, ID is a unique identifier (e.g., wallet address) of the party for whom a verdict is provided (e.g., a client), and  $o$  is a counter that determines how many times a verdict for the same ID holder has been generated in the past. ID and  $o$  are used to ensure that each  $r_j$  will be different for each invocation of PVE although the same key  $\bar{k}_0$  is used.

**FVD( $n, \bar{w}$ )  $\rightarrow v$**

- **Input.**  $n$ : the total number of auditors, and  $\bar{w} = [\bar{w}_1, \dots, \bar{w}_n]$ : a vector of all auditors’ encoded verdicts.
- **Output.**  $v$ : final verdict.

A third-party  $\mathcal{I}$  takes the following steps.

1) combines all auditors’ encoded verdicts,  $\bar{w}_j \in \bar{w}$ , as follows.  $c = \bigoplus_{j=1}^n \bar{w}_j$

2) sets the final verdict  $v$  depending on the content of  $c$ . Specifically,

$$v = \begin{cases} 0, & \text{if } c = 0 \\ 1, & \text{otherwise} \end{cases}$$

3) outputs  $v$ .

Figure 4: Final Verdict Decoding (FVD) Protocol.

## 7.5. The PwDR Protocol

In this section, we present the PwDR protocol in detail.

1) **Generating  $\mathcal{G}$ ’s and  $\mathcal{D}_j$ ’s Parameters:**  $\text{keyGen}(1^\lambda) \rightarrow (\bar{sk}, pk)$ . Parties  $\mathcal{G}$  and (only)  $\mathcal{D}_j$  take steps 1a and 1b respectively.

a) calls  $\text{Sig.keyGen}(1^\lambda) \rightarrow (sk_{\mathcal{G}}, pk_{\mathcal{G}})$  to generate secret key  $sk_{\mathcal{G}}$  and public key  $pk_{\mathcal{G}}$ . It publishes  $pk_{\mathcal{G}}$ .

b) calls  $\text{keyGen}(1^\lambda) \rightarrow (\bar{sk}_{\mathcal{D}}, pk_{\mathcal{D}})$  to generate decrypting secret key  $\bar{sk}_{\mathcal{D}}$  and encrypting public key  $pk_{\mathcal{D}}$ . It also generates a key  $\bar{k}_0$  for PRF, i.e.,  $\bar{k}_0 \xleftarrow{\$} \{0, 1\}^\lambda$ . It sets  $pk_{\mathcal{D}} = \bar{pk}_{\mathcal{D}}$  and  $sk_{\mathcal{D}} := (\bar{sk}_{\mathcal{D}}, \bar{k}_0)$ . It publishes  $pk_{\mathcal{D}}$  and sends  $sk_{\mathcal{D}}$  to the rest of the auditors.

1. It is similar to the idea used in the XOR-based secret sharing [40].

Let  $sk := (sk_G, sk_D)$  and  $pk := (pk_G, pk_D)$ . Note, this phase occurs only once for all customers.

2) Bank-side Initiation:  $\text{bankInit}(1^\lambda) \rightarrow (T, pp, l)$ .

Bank  $\mathcal{B}$  takes the following steps.

- picks secret keys  $\bar{k}_1$  and  $\bar{k}_2$  for the symmetric key encryption scheme. It sets two private statements as  $\pi_1 = \bar{k}_1$  and  $\pi_2 = \bar{k}_2$ .
- calls  $\text{SAP.init}(1^\lambda, \text{adr}_B, \text{adr}_C, \pi_i) \rightarrow (r_i, g_i, \text{adr}_{\text{SAP}})$  to initiate agreements on statements  $\pi_i \in \{\pi_1, \pi_2\}$  with  $\mathcal{C}$ . Let  $T_i := (\pi_i, g_i)$  and  $T := (T_1, T_2)$ , where  $\pi_i := (\pi_i, r_i)$  is the opening of  $g_i$ . It also sets parameter  $\Delta$  as a time window between two specific time points, i.e.,  $\Delta = t_i - t_{i-1}$ . Briefly, it is used to impose an upper bound on a message delay.
- sends  $\tilde{\pi} := (\tilde{\pi}_1, \tilde{\pi}_2)$  to  $\mathcal{C}$  and sends public parameter  $pp := (\text{adr}_{\text{SAP}}, \Delta)$  to smart contract  $\mathcal{S}$ .

3) Customer-side Initiation:

$\text{customerInit}(1^\lambda, T, pp) \rightarrow a$ .

Customer  $\mathcal{C}$  takes the following steps.

- calls  $\text{SAP.agree}(\pi_i, r_i, g_i, \text{adr}_B, \text{adr}_{\text{SAP}}) \rightarrow (g'_i, b_i)$ , to locally check the correctness of parameters in  $T_i \in T$  and (if accepted) to agree on these parameters, where  $(\pi_i, r_i) \in \tilde{\pi}_i \in T_i$  and  $1 \leq i \leq 2$ . Note, if both  $\mathcal{B}$  and  $\mathcal{C}$  are honest, then  $g_i = g'_i$ . It also checks  $\Delta$  in  $\mathcal{S}$ , e.g., to see if it is sufficiently large.
- if the above checks fail, it sets  $a = 0$  and aborts. Otherwise, it sets  $a = 1$ . It sends  $a$  to  $\mathcal{S}$ .

4) Generating Update Request:

$\text{genUpdateRequest}(T, f, l) \rightarrow \hat{m}_1^{(C)}$ .

Customer  $\mathcal{C}$  takes the following steps.

- sets its request parameter  $m_1^{(C)}$  as below.
  - if it wants to set up a new payee, then it sets  $m_1^{(C)} := (\phi, f)$ , where  $f$  is the new payee's detail.
  - if it wants to amend the existing payee's detail, it sets  $m_1^{(C)} := (i, f)$ , where  $i$  is an index of the element in  $l$  that should change to  $f$ .
- at time  $t_0$ , sends to  $\mathcal{S}$  the encryption of  $m_1^{(C)}$ , i.e.,  $\hat{m}_1^{(C)} = \text{Enc}(\bar{k}_1, m_1^{(C)})$ .

5) Inserting New Payee:  $\text{insertNewPayee}(\hat{m}_1^{(C)}, l) \rightarrow \hat{l}$ .

Smart contract  $\mathcal{S}$  takes the following steps.

- if  $\hat{m}_1^{(C)}$  is not empty, it appends  $\hat{m}_1^{(C)}$  to the payee list  $\hat{l}$ , resulting in an updated list,  $\hat{l}$ .
- if  $\hat{m}_1^{(C)}$  is empty, it does nothing.

6) Generating Warning:  $\text{genWarning}(T, \hat{l}, aux) \rightarrow \hat{m}_1^{(B)}$ .

Bank  $\mathcal{B}$  takes the following steps.

- checks if the most recent list  $\hat{l}$  is not empty. If it is empty, it halts. Else, it proceeds to the next step.
- decrypts each element of  $\hat{l}$  and checks its correctness, e.g., checks whether each element meets its internal policy stated in  $aux$ . If the check passes, it sets  $m_1^{(B)} = \text{"pass"}$ . Otherwise, it sets  $m_1^{(B)} = \text{"warning"}$ , where the warning is a string that contains a warning's detail concatenated with the string "warning".
- at time  $t_1$ , sends to  $\mathcal{S}$  the encryption of  $m_1^{(B)}$ , i.e.,  $\hat{m}_1^{(B)} = \text{Enc}(\bar{k}_1, m_1^{(B)})$ .

7) Generating Payment Request:

$\text{genPaymentRequest}(T, in_f, \hat{l}, \hat{m}_1^{(B)}) \rightarrow \hat{m}_2^{(C)}$ .

Customer  $\mathcal{C}$  takes the following steps.

- at time  $t_2$ , decrypts  $\hat{l}$  and  $\hat{m}_1^{(B)}$ . Depending on the warning, it sets a payment request  $m_2^{(C)}$  to  $\phi$  or  $in_f$ , where  $in_f$  contains the payment's detail, e.g., the payee's detail in  $\hat{l}$  and amount it wants to send.
- at time  $t_3$ , sends to  $\mathcal{S}$  the encryption of  $m_2^{(C)}$ , i.e.,  $\hat{m}_2^{(C)} = \text{Enc}(\bar{k}_1, m_2^{(C)})$ .

8) Making Payment:  $\text{makePayment}(T, \hat{m}_2^{(C)}) \rightarrow \hat{m}_2^{(B)}$ .

Bank  $\mathcal{B}$  takes the following steps.

- at time  $t_4$ , decrypts  $\hat{m}_2^{(C)}$ , i.e.,  $m_2^{(C)} = \text{Dec}(\bar{k}_1, \hat{m}_2^{(C)})$ .
- at time  $t_5$ , checks the content of  $m_2^{(C)}$ . If  $m_2^{(C)}$  is non-empty, i.e.,  $m_2^{(C)} = in_f$ , it checks if the payee's detail in  $in_f$  has already been checked and the payment's amount does not exceed the customer's credit. If the checks pass, it runs the off-chain payment algorithm,  $\text{pay}(in_f)$ . In this case, it sets  $m_2^{(B)} = \text{"paid"}$ . Otherwise (i.e., if  $m_2^{(C)} = \phi$  or neither checks pass), it sets  $m_2^{(B)} = \phi$ . It sends to  $\mathcal{S}$  the encryption of  $m_2^{(B)}$ , i.e.,  $\hat{m}_2^{(B)} = \text{Enc}(\bar{k}_1, m_2^{(B)})$ .

9) Generating Complaint:

$\text{genComplaint}(\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f) \rightarrow (\hat{z}, \hat{\pi})$ .

Customer  $\mathcal{C}$  takes the following steps.

- decrypts  $\hat{m}_1^{(B)}$  and  $\hat{m}_2^{(B)}$ ; this results in  $m_1^{(B)}$  and  $m_2^{(B)}$  respectively. Depending on the content of the decrypted values, it sets its complaint's parameters  $z := (z_1, z_2, z_3)$  as follows.
  - if  $\mathcal{C}$  wants to make one of the two below statements, it sets  $z_1 = \text{"challenge message"}$ .
    - the pass message (in  $m_1^{(B)}$ ) should have been a warning.
    - $\mathcal{B}$  did not provide any message and if  $\mathcal{B}$  provided a warning, the fraud would have been prevented.
  - if  $\mathcal{C}$  wants to challenge the effectiveness of the warning (in  $m_1^{(B)}$ ), it sets  $z_2 = m || \text{sig} || pk_G || \text{"challenge warning"}$ , where  $m$  is a piece of evidence,  $\text{sig} \in aux_f$  is the evidence's certificate (obtained from  $\mathcal{G}$ ), and  $pk_G \in pk$ .
  - if  $\mathcal{C}$  wants to complain about the payment's inconsistency, it sets  $z_3 = \text{"challenge payment"}$ ; else, it sets  $z_3 = \phi$ .

- at time  $t_6$ , sends  $\hat{z} = \text{Enc}(\bar{k}_1, z)$  and  $\hat{\pi} = \text{Enc}(pk_D, \tilde{\pi})$  to  $\mathcal{S}$ .

10) Verifying Complaint:  $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, \frac{sk_D}{sk_D}, aux, pp) \rightarrow \hat{w}_j$ .

Every  $\mathcal{D}_j \in \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$  acts as follows.

- at time  $t_7$ , decrypts  $\hat{\pi}$ , i.e.,  $\tilde{\pi} = \text{Dec}(\tilde{sk}_D, \hat{\pi})$ , where  $\tilde{sk}_D \in sk_D$ .
- checks the validity of  $(\tilde{\pi}_1, \tilde{\pi}_2)$  in  $\tilde{\pi}$  by locally running the SAP's verification, i.e.,  $\text{SAP.verify}(\cdot)$ , for each  $\tilde{\pi}_i$ . It returns  $s$ . If  $s = 0$ , it halts. If  $s = 1$  for both  $\tilde{\pi}_1$  and  $\tilde{\pi}_2$ , it proceeds to the next step.
- decrypts  $\hat{m} = [\hat{m}_1^{(C)}, \hat{m}_2^{(C)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ , using  $\text{Dec}(\bar{k}_1, \cdot)$ , where  $\bar{k}_1 \in \tilde{\pi}_1$ . Let  $[m_1^{(C)}, m_2^{(C)}, m_1^{(B)}, m_2^{(B)}]$  be the result.
- checks whether  $\mathcal{C}$  made an update request to its payee's list. To do so, it checks if  $m_1^{(C)}$  is non-empty and (its encryption) was registered by  $\mathcal{C}$  in  $\mathcal{S}$ . Also, it checks whether  $\mathcal{C}$  made a payment request, by checking if  $m_2^{(C)}$  is non-empty and (its encryption) was registered by  $\mathcal{C}$  in  $\mathcal{S}$  at time  $t_3$ . If either check

fails, it halts.

e) decrypts  $\hat{z}$  and  $\hat{l}$  using  $\text{Dec}(\bar{k}_1, \cdot)$ , where  $\bar{k}_1 \in \hat{\pi}_1$ . Let  $z := (z_1, z_2, z_3)$  and  $l$  be the result.

f) sets its verdicts according to the content of  $z := (z_1, z_2, z_3)$ , as follows.

- if “challenge message”  $\notin z_1$ , it sets  $w_{1,j} = 0$ . Otherwise, it runs  $\text{verStat}(\text{add}_S, m_1^{(B)}, l, \Delta, \text{aux}) \rightarrow w_{1,j}$ , to determine if a warning (in  $m_1^{(B)}$ ) should have been given (instead of the pass or no message).
- if “challenge warning”  $\notin z_2$ , it sets  $w_{2,j} = w_{3,j} = 0$ . Otherwise, it runs  $\text{checkWarning}(\text{add}_S, z_2, m_1^{(B)}, \text{aux}') \rightarrow (w_{2,j}, w_{3,j})$ , to determine the effectiveness of the warning (in  $m_1^{(B)}$ ).
- if “challenge payment”  $\in z_3$ , it checks if the payment was made. If it passes, it sets  $w_{4,j} = 1$ . If it fails, it sets  $w_{4,j} = 0$ . If “challenge payment”  $\notin z_3$ , it checks if “paid”  $\in m_2^{(B)}$ . If it passes, it sets  $w_{4,j} = 1$ . Else, it sets  $w_{4,j} = 0$ .

g) encodes its verdicts  $(w_{1,j}, w_{2,j}, w_{3,j}, w_{4,j})$ :

- i) locally maintains a counter,  $o_{\text{adv}_C}$ , for each  $C$ . It sets its initial value to 0. It skips this step, if the counter has already been set.
- ii) calls  $\text{PVE}(\cdot)$  to encode each verdict. In particular, it performs as follows.  $\forall i, 1 \leq i \leq 4$ :
  - calls  $\text{PVE}(\bar{k}_0, \text{adv}_C, w_{i,j}, o_{\text{adv}_C}, n, j) \rightarrow \bar{w}_{i,j}$ , where  $\bar{k}_0 \in sk_D$ .
  - sets  $o_{\text{adv}_C} = o_{\text{adv}_C} + 1$ .

By the end of this step, a vector  $\bar{w}_j$  of four encoded verdicts is computed, i.e.,  $\bar{w}_j = [\bar{w}_{1,j}, \dots, \bar{w}_{4,j}]$ .

iii) uses  $\bar{k}_2 \in \hat{\pi}_2$  to further encode/encrypt  $\text{PVE}(\cdot)$ 's outputs as follows.  $\hat{w}_j = \text{Enc}(\bar{k}_2, \bar{w}_j)$ .

h) at time  $t_8$ , sends to  $S$  the encrypted vector,  $\hat{w}_j$ .

11) *Resolving Dispute*:  $\text{resDispute}(T_2, \hat{w}, pp) \rightarrow v$ .

Party  $\mathcal{DR}$  takes the below steps at time  $t_9$ , if it is invoked by  $C$  or  $S$  which sends  $\hat{\pi}_2 \in T_2$  to it.

- a) checks  $\hat{\pi}_2$ 's validity by locally running the SAP's verification, i.e.,  $\text{SAP.verify}(\cdot)$ , that returns  $s$ . If  $s = 0$ , it halts.
- b) computes the final verdicts, as below.
  - i) uses  $\bar{k}_2 \in \hat{\pi}_2$  to decrypt the auditors' encoded verdicts, as follows.  $\forall j, 1 \leq j \leq n : \bar{w}_j = \text{Dec}(\bar{k}_2, \hat{w}_j)$ , where  $\hat{w}_j \in \hat{w}$ .
  - ii) constructs four vectors,  $[u_1, \dots, u_4]$ , and sets each vector  $u_i$  as follows.  $\forall i, 1 \leq i \leq 4 : u_i = [\bar{w}_{i,1}, \dots, \bar{w}_{i,n}]$ , where  $\bar{w}_{i,j} \in \bar{w}_j$ .
  - iii) calls  $\text{FVD}(\cdot)$  to extract each final verdict, as follows.  $\forall i, 1 \leq i \leq 4 : \text{calls FVD}(n, u_i) \rightarrow v_i$ .
- c) outputs  $v = [v_1, \dots, v_4]$ .

Customer  $C$  must be reimbursed if the final verdict is that (i) the “pass” message or missing message should have been a warning or (ii) the warning was ineffective and the provided evidence was not invalid, and (iii) the payment has been made. To state it formally, the following relation must hold:

$$\left( \underbrace{(v_1 = 1)}_{(i)} \vee \underbrace{(v_2 = 1 \wedge v_3 = 1)}_{(ii)} \right) \wedge \left( \underbrace{v_4 = 1}_{(iii)} \right)$$

In the above PwDR protocol, even  $C$  and  $B$  that know

the decryption secret keys,  $(\bar{k}_1, \bar{k}_2)$ , cannot link a certain verdict to an auditor, because: (a) they do not know the masking random values used by auditors to mask each verdict and (b) the final verdicts  $(v_1, \dots, v_4)$  reveal nothing about the number of 1 or 0 verdicts, except when all auditors vote 0. In the PwDR, we used PVE and FVD only because they are highly efficient. However, it is easy to replace them with GPVE and GFVD.

We also highlight that our protocol does not require the bank to commit any funds to the smart contract, which keeps it consistent with the traditional banking setting. The final (legal) verdict would suffice to enforce the bank to reimburse victims. Furthermore, in the real world, during a payment journey, a customer may receive various warning messages depending on the details it provides, its transaction history, and the checks a bank conducts, e.g., “Confirmation of Payee” [14]. Thus, we have included warning messages in our protocol to match the real-world banking setting.

Below, we present the security theorem of the PwDR protocol and present its formal proof in Appendix J.

**Theorem 1.** *The above PwDR protocol is secure, with regard to Definition 12, if the digital signature is existentially unforgeable under chosen message attacks, SAP, and the verdict encoding-decoding protocols (i.e., PVE and FVD) are secure, the encryption schemes are semantically secure, the blockchain is immutable, and the correctness of PVE and FVD holds.*

Note that since the smart contract in our PwDR protocol merely acts as an immutable bulletin board, one may replace it with any other efficient tamper-evident logging mechanism, e.g., [15], [30], [41].

## 8. Evaluation

In this section, we analyse the PwDR protocol's (computation and communication) complexity, its concrete runtime, and transaction latency. Tables 2 and 3 summarize the asymptotic and performance analysis respectively.

### 8.1. Computation Complexity

We first analyse  $C$ 's cost. In Phase 3,  $C$  invokes a hash function twice to check the correctness of the private statements' parameters. In Phase 4, it invokes the symmetric encryption once to encrypt its update request. In Phase 7, it invokes the symmetric encryption twice to decrypt  $B$ 's warning message and to encrypt its payment request. In Phase 9, it runs the symmetric encryption three times to decrypt  $B$ 's warning and payment messages and to encrypt its complaint. In the same phase, it invokes asymmetric encryption once to encrypt the private statements' opening. Therefore,  $C$ 's complexity is  $O(1)$ . Next, we analyse  $B$ 's cost. In Phase 2, it invokes the hash function twice to commit to two statements. In Phase 6, it calls the symmetric key encryption once to encrypt its outgoing warning message. In Phase 8, it also invokes the symmetric key encryption once to encrypt the outgoing payment message. Thus,  $B$ 's complexity is  $O(1)$  too. Next, we analyse each auditor's cost. In Phase 10, each  $\mathcal{D}_j$  invokes the asymmetric key encryption once to decrypt the private statements' openings. It also invokes the hash function twice to verify the openings. It invokes the symmetric key

TABLE 2: The PwDR’s asymptotic cost. In the table,  $n$  is the number of auditors and  $e$  is the threshold.

Party	Setting		Computation Cost	Communication Cost
	$e = 1$	$e > 1$		
Customer $\mathcal{C}$	✓	✓	$O(1)$	$O(1)$
Bank $\mathcal{B}$	✓	✓	$O(1)$	$O(1)$
Auditor $\mathcal{D}_1, \dots, \mathcal{D}_{n-1}$	✓	✓	$O(1)$	$O(1)$
Auditor $\mathcal{D}_n$	✓		$O(n)$	$O(1)$
		✓	$O(\sum_{i=e}^n \frac{n!}{i!(n-i)!})$	$O(\sum_{i=e}^n \frac{n!}{i!(n-i)!})$
Dispute resolver $\mathcal{DR}$	✓	✓	$O(n)$	$O(1)$

TABLE 3: Parties’ run-time (in ms) in verdict encoding-decoding protocols.  $n$ : the number of auditors and  $e$ : the threshold.

Party	$n = 6$		$n = 8$		$n = 10$		$n = 12$	
	$e = 1$	$e = 4$	$e = 1$	$e = 5$	$e = 1$	$e = 6$	$e = 1$	$e = 7$
Auditor $\mathcal{D}_n$	0.019	0.220	0.033	0.661	0.035	2.87	0.052	10.15
Dispute resolver $\mathcal{DR}$	0.001	0.015	0.001	0.016	0.001	0.069	0.003	0.09

encryption six times to decrypt  $\mathcal{C}$ ’s and  $\mathcal{B}$ ’s messages that were posted on  $\mathcal{S}$  (this includes  $\mathcal{C}$ ’s complaint). Recall, in the same phase, each auditor encodes its verdict using a verdict encoding protocol. Now, we evaluate the verdict encoding complexity of each auditor for two cases: (a)  $e = 1$  and (b)  $e \in (1, n]$ . Note, in the former case the PVE is invoked while in the latter GPVE is invoked. In case (a), every auditor  $\mathcal{D}_j$ , except  $\mathcal{D}_n$ , invokes the pseudorandom function once to encode its verdict. However, auditor  $\mathcal{D}_n$  invokes the pseudorandom function  $n - 1$  times and XORs the function’s outputs with each other. Thus, in case (a), auditor  $\mathcal{D}_n$ ’s complexity is  $O(n)$  while the rest of auditors’ complexity is  $O(1)$ . In case (b), every auditor  $\mathcal{D}_j$ , except  $\mathcal{D}_n$ , invokes the pseudorandom function twice to encode its verdict. But, auditor  $\mathcal{D}_n$  invokes the pseudorandom function  $n - 1$  times and XORs the function’s outputs with each other. It invokes the pseudorandom function  $n$  times to generate all auditors’ representations of verdict

1. It computes all  $y = \sum_{i=e}^n \frac{n!}{i!(n-i)!}$  combinations of the representations that meet the threshold which involves  $O(y)$  XORs. It inserts  $y$  elements into a Bloom filter that requires  $O(y)$  hash function evaluations. So, in case (b), auditor  $\mathcal{D}_n$ ’s complexity is  $O(y)$  while the rest of the auditors’ complexity is  $O(1)$ . To conclude, in Phase 10, auditor  $\mathcal{D}_n$ ’s complexity is either  $O(n)$  or  $O(y)$ , while the rest of the auditors’ complexity is  $O(1)$ . Now, we analyse  $\mathcal{DR}$ ’s cost in Phase 11. It invokes the hash function once to check the private statement’s correctness. It also performs  $O(n)$  symmetric key decryption to decrypt auditors’ encoded verdicts. Now, we evaluate the verdict decoding complexity of  $\mathcal{DR}$  for two cases: (a)  $e = 1$  and (b)  $e \in (1, n]$ . In the former case (in which FVD is invoked), it performs  $O(n)$  XOR to combine all verdicts. Its complexity is  $O(n)$  in the latter case (in which GFVD is invoked), with the difference that it also invokes the Bloom filter’s hash functions, to make a membership query to the Bloom filter. Thus,  $\mathcal{DR}$ ’s complexity is  $O(n)$ .

## 8.2. Communication Cost

Now, we analyse the communication cost of the PwDR protocol. Briefly,  $\mathcal{C}$ ’s complexity is  $O(1)$  as in total it sends only six messages to other parties. Similarly,  $\mathcal{B}$ ’s complexity is  $O(1)$  as its total number of outgoing mes-

sages is only nine. Each auditor  $\mathcal{D}_j$  sends only four messages to the smart contract, so its complexity is  $O(1)$ . However, if GFVD is invoked, then auditor  $\mathcal{D}_n$  needs to send also a Bloom filter that costs it  $O(y)$ . Moreover,  $\mathcal{DR}$ ’s complexity is  $O(1)$ , as its outgoing messages include only four binary values.

## 8.3. Concrete Performance Analysis

In this section, we study the protocol’s performance. As we saw in the previous section, the customer’s and bank’s complexity is very low and constant; however, one of the auditors, i.e., auditor  $\mathcal{D}_n$ , and the dispute resolver have non-constant complexities. These non-constant overheads were mainly imposed by the verdict inducing-decoding protocols. Therefore, to study these parties’ runtime in the PwDR, we implemented both variants of the verdict encoding-decoding protocols (that were presented in Section 7.4). They were implemented in C++, see [5], [6] for the source code. To conduct the experiment, we used a MacBook Pro laptop with quad-core Intel Core i5, 2 GHz CPU, and 16 GB RAM. We ran the experiment on average 100 times. The prototype implementation uses the “Cryptopp” library<sup>2</sup> for cryptographic primitives, the “GMP” library<sup>3</sup> for arbitrary precision arithmetic, and the “Bloom Filter” library<sup>4</sup>. In the experiment, we set the false-positive rate in a Bloom filter to  $2^{-40}$  and the finite field size to 128 bits. We used AES to implement PRF. Table 3 (in Section 8) provides the runtime of  $\mathcal{D}_n$  and  $\mathcal{DR}$  for various numbers of auditors in two cases; namely, when the threshold is 1 and when it is greater than 1. In the former case, we used the PVE and FVD protocols. In the latter case, we used the GPVE and GFVD ones.

As Table 3 depicts, the runtime of  $\mathcal{D}_n$  increases gradually from 0.019 to 10.15 milliseconds when the number of auditors grows from  $n = 6$  to  $n = 12$ . In contrast, the runtime of  $\mathcal{DR}$  grows slower; it increases from 0.001 to 0.09 milliseconds when the number of auditors increases. Nevertheless, the overall cost is very low. Specifically, the highest runtime is only about 10 milliseconds which belongs to  $\mathcal{D}_n$  when  $n = 12$  and  $e = 7$ . It is also evident

2. <https://www.cryptopp.com>

3. <https://gmplib.org>

4. <http://www.partow.net/programming/bloomfilter/index.html>

that the parties' runtime in the PVE and FVD protocols is much lower than their runtime in the GPVE and GFVD ones. To compare the parties' runtime, we also fixed the threshold to 6 (in GPVE and GFVD protocols) and ran the experiment for different values of  $n$ . Figure 5 summarises the result. As this figure indicates, the runtime of  $\mathcal{D}_n$  and  $\mathcal{DR}$  almost linearly grows when the number of auditors increases. Moreover,  $\mathcal{D}_n$  has a higher runtime than  $\mathcal{DR}$  has, and its runtime growth is faster than that of  $\mathcal{DR}$ .

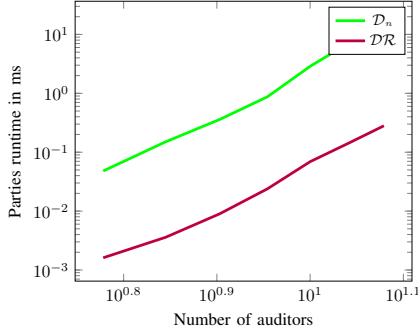


Figure 5: Parties' runtime in the PwDR.

## 8.4. Transaction Latency

The transaction latency imposed by the blockchain to the PwDR protocol depends on the type of consensus protocol used. For instance, on average it takes between 12 seconds in Ethereum to mine a block [19]; after this block is propagated to the network, to have adequate confidence that the block will remain in the chain, one may need to wait until at least 6 blocks are added after that block, which would take about 72 additional seconds. However, such a delay is lower in Byzantine Fault Tolerance (BFT) Hyperledger Fabric blockchain as it does not involve any mining and a consensus can be reached faster, e.g., about 35 seconds when 20 nodes are involved, [25].

## 9. Future Research

In this section, we highlight a set of future research directions in the context of APP frauds.

### 9.1. Improving Warnings' Effectiveness

As we stated in Section 5.2, one of the determining factors in the process of allocating liability to APP fraud victims is if they follow warnings. However, there exists no publicly available study on the effectiveness of warnings in the context of APP frauds. There exists a comprehensive research line in determining the effectiveness of warnings in general, e.g., in [12], [21], [29]. Nevertheless, in the context of APP fraud, there is a vital unique factor that can directly influence a warning's effectiveness. The factor is the ability of fraudsters to interact directly with their victims. This lets a fraudster actively try to negate the effectiveness of a bank's warning and persuade the warning recipient to ignore the warning (and make payment). Such a factor was not taken into account in the traditional study of warnings. The current high rate of APP fraud occurrence suggests that there exists a huge room for improving warnings' effectiveness. Thus, future research can identify key factors that can improve the effectiveness of warnings in this context.

### 9.2. Protecting APP Fraud Victims of Alternative Payment Platforms

To date, there is no official report on the occurrence of APP frauds on any payment platforms (e.g., cryptocurrencies) other than regular banking. This can be due to a lack of an oversight organisation which collects fraud-related data or due to a low rate of such fraud taking place on these platforms because these platforms are not popular enough. However, with the increase in the popularity of alternative payment platforms (including Central Bank Digital Currency), it is likely that APP fraudsters will target these platforms' users. Hence, another interesting future research direction would be to design secure dispute resolution protocols to protect APP fraud victims on these platforms as well.

### 9.3. Studying Users' Compliance with the CRM Code's Guidelines

Currently, customers are expected to comply with the CRM code's guidelines. Users' compliance with these guidelines could help lower the rate of APP fraud occurrence. If victims fail to comply with such guidelines, then banks can claim that customers' have been negligent which ultimately could cost the victims. Such guidelines would be effective when they are known and followed by customers. Recently, Van Der Zee [50] has conducted a study to find out whether customers of the "Dutch Banking Association" are aware of the bank's digital payments guidelines and if so, whether they comply with these guidelines. But, there exists no systematic study to investigate whether customers are aware of and comply with the CRM code's guidelines. Therefore, another research direction is to fill the above void.

### 9.4. Security Against Exploitative Victims

Having in place a transparent deterministic procedure (e.g., the PwDR protocol) for evaluating victims' requests for reimbursement could potentially create opportunities for exploitation. In particular, an honest victim of an APP fraud that had been reimbursed in the past due to the payment system's vulnerability (e.g., an ineffective warning) may be tempted to exploit the same known vulnerability multiple times. Hence, future research can investigate how to secure the online banking system against such exploitative victims.

## 10. Conclusion

In this work, to facilitate APP frauds victims' reimbursement, we proposed the notion of payment with dispute resolution. We identified the vital properties that such a notion should possess and formally defined them. We also proposed a candidate construction, PwDR, and proved its security. The PwDR not only offers transparency and accountability but also acts as a data hub providing sufficient information that could help regulators examine whether the reimbursement regulations have been applied correctly and consistently among financial institutions. We also studied the PwDR's cost via asymptotic and concrete runtime evaluation. Our cost analysis indicated that the construction is indeed efficient.

## References

- [1] A P20 Fraud and Criminal Transactions Working Group Paper. Best practice approaches for combating payee scams. 2021. <https://static1.squarespace.com/static/5efcc6dae323db37b4d01d19/t/60958e7453a1e0728b445470/1620414069082/P20+Report+-+Best+Practice+Approaches+For+Combating+Payee+Scams.pdf>.
- [2] Aydin Abadi, Steven J. Murdoch, and Thomas Zacharias. Recurring contingent payment for proofs of retrievability. *IACR Cryptol. ePrint Arch.*, 2021.
- [3] Ross Anderson, Chris Barton, Rainer Bölme, Richard Clayton, Carlos Ganán, Tom Grasso, Michael Levi, Tyler Moore, and Marie Vasek. Measuring the changing cost of cybercrime. 2019.
- [4] Ross Anderson et al. Closing the phishing hole—fraud, risk and nonbanks. In *Federal Reserve Bank of Kansas City—Payment System Research Conferences*, pages 41–56, 2007.
- [5] Anonymous. Variant 1: Efficient verdict encoding-decoding protocol, 2021. <https://github.com/pwdrprotocol/PwDR/blob/main/encoding-decoding.cpp>.
- [6] Anonymous. Variant 2: Generic verdict encoding-decoding protocol, 2021. <https://github.com/pwdrprotocol/PwDR/blob/main/generic-encoding-decoding.cpp>.
- [7] Financial Conduct Authority. FCA glossary, 2021. <https://www.handbook.fca.org.uk/handbook/glossary/G3566a.html>.
- [8] Ingolf Becker, Alice Hutchings, Ruba Abu-Salma, Ross J. Anderson, Nicholas Bohm, Steven J. Murdoch, M. Angela Sasse, and Gianluca Stringhini. International comparison of bank fraud reimbursement: customer perceptions and contractual terms. *J. Cybersecur.*, 2017.
- [9] Burton H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Commun.*, 1970.
- [10] Nicholas Bohm, Ian Brown, and Brian Gladman. Electronic commerce: Who carries the risk of fraud? *J. Inf. Law Technol.*, 2000, 2000.
- [11] Prosenjit Bose, Hua Guo, Evangelos Kranakis, Anil Maheshwari, Pat Morin, Jason Morrison, Michiel H. M. Smid, and Yihui Tang. On the false-positive rate of bloom filters. *Inf. Process. Lett.*, 2008.
- [12] Bonnie Brinton Anderson, Anthony Vance, C Brock Kirwan, David Eargle, and Jeffrey L Jenkins. How users perceive and respond to security messages: a neurois research agenda and empirical study. *European Journal of Information Systems*, 25(4):364–390, 2016.
- [13] Michael Buchwald. Smart contract dispute resolution: the inescapable flaws of blockchain-based arbitration. *U. Pa. L. Rev.*, 2019.
- [14] Confirmation of Payee Team. Confirmation of payee- response to consultation cp20/1 and decision on varying specific direction 10. 2020. <https://www.psr.org.uk/media/qrb03jm/psr-ps20-1-variation-of-specific-direction-10-february-2020.pdf>.
- [15] Scott A. Crosby and Dan S. Wallach. Efficient data structures for tamper-evident logging. In Fabian Monrose, editor, *18th USENIX Security Symposium, Montreal, Canada, August 10-14, 2009, Proceedings*, pages 317–334. USENIX Association, 2009.
- [16] Changyu Dong, Yilei Wang, Amjad Aldweesh, Patrick McCorry, and Aad van Moorsel. Betrayal, distrust, and rationality: Smart counter-collusion contracts for verifiable cloud computing. In *CCS*, 2017.
- [17] Stefan Dziembowski, Lisa Ekey, and Sebastian Faust. Fairswap: How to fairly exchange digital goods. In *CCS*, 2018.
- [18] Lisa Ekey, Sebastian Faust, and Benjamin Schlosser. Optiswap: Fast optimistic fair exchange. In *ASIA CCS*, 2020.
- [19] Ittay Eyal, Adem Efe Gencer, Emin Gün Sirer, and Robbert van Renesse. Bitcoin-ng: A scalable blockchain protocol. In Katerina J. Argyraki and Rebecca Isaacs, editors, *13th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2016, Santa Clara, CA, USA, March 16-18, 2016*, pages 45–59. USENIX Association, 2016.
- [20] Federal Bureau of Investigation (FBI). Internet crime report. 2020. [https://www.ic3.gov/Media/PDF/AnnualReport/2020\\_IC3Report.pdf](https://www.ic3.gov/Media/PDF/AnnualReport/2020_IC3Report.pdf).
- [21] Adrienne Porter Felt, Robert W Reeder, Hazim Almuhammedi, and Sunny Consolvo. Experimenting at scale with google chrome’s SSL warning. In *CHI*, 2014.
- [22] Adam French. Which? makes scams super-complaint-banks must protect those tricked into a bank transfer. 2016. <https://www.which.co.uk/news/2016/09/which-makes-scams-super-complaint-453196/>.
- [23] Juan A. Garay, Aggelos Kiayias, and Nikos Leonardos. The bitcoin backbone protocol: Analysis and applications. In *EUROCRYPT*, 2015.
- [24] Matthew Green and Ian Miers. Bolt: Anonymous payment channels for decentralized currencies. In *CCS*, 2017.
- [25] Hyperledger Foundation. Hyperledger blockchain performance metrics, 2018. [https://www.hyperledger.org/wp-content/uploads/2018/10/HL\\_Whitepaper\\_Metrics\\_PDF\\_V1.01.pdf](https://www.hyperledger.org/wp-content/uploads/2018/10/HL_Whitepaper_Metrics_PDF_V1.01.pdf).
- [26] Jonathan Katz and Yehuda Lindell. *Introduction to Modern Cryptography, Second Edition*. CRC Press, 2014.
- [27] Marte Eidsand Kjørven. Who pays when things go wrong? online financial fraud and consumer protection in scandinavia and europe. *European Business Law Review*, 2020.
- [28] Ralf Küsters, Julian Liedtke, Johannes Müller, Daniel Rausch, and Andreas Vogt. Ordinos: A verifiable tally-hiding e-voting system. In *EuroS&P*, 2020.
- [29] Kenneth R Laughery and Michael S Wogalter. Designing effective warnings. *Reviews of human factors and ergonomics*, 2006.
- [30] Ben Laurie, Adam Langley, and Emilia Käsper. Certificate transparency. *RFC*, 6962:1–27, 2013.
- [31] Lending Standards Board. Contingent reimbursement model code for authorised push payment scams. 2021. <https://www.lendingstandardsboard.org.uk/wp-content/uploads/2021/04/CRM-Code-LSB-Final-April-2021.pdf>.
- [32] David McIlroy and Ruhi Sethi-Smith. Prospects for bankers’ liability for authorised push payment fraud. *Butterworths Journal of International Banking and Financial Law*, 2021. <https://www.forumchambers.com/wp-content/uploads/2021/03/Article-2-Smith.1-1.pdf>.
- [33] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system. Technical report, 2019.
- [34] US House Committee on Financial Services. Discussion draft: To amend the electronic fund transfer act to treat fraudulently induced electronic fund transfers in the same manner as unauthorized electronic fund transfer, 2022. <https://docs.house.gov/meetings/BA/BA00/20220428/114690/BILLS-117pih-ProtectingConsumersFromPaym-U1.pdf>.
- [35] Pietro Ortolani. Self-enforcing online dispute resolution: lessons from bitcoin. *Oxford Journal of Legal Studies*, 2016.
- [36] Pietro Ortolani. The impact of blockchain technologies and smart contracts on dispute resolution: arbitration and court litigation at the crossroads. *Uniform law review*, 2019.
- [37] Payment Systems Regulator. What are authorised push payment scams, 2022. <https://www.psr.org.uk/our-work/app-scams>.
- [38] Torben P. Pedersen. Non-interactive and information-theoretic secure verifiable secret sharing. In *CRYPTO*, 1991.
- [39] Joseph Poon and Thaddeus Dryja. The bitcoin lightning network: Scalable off-chain instant payments. Technical report, 2016. <https://lightning.network/lightning-network-paper.pdf>.
- [40] Bruce Schneier. *Applied cryptography - protocols, algorithms, and source code in C, 2nd Edition*. Wiley, 1996.
- [41] Enrique Soriano-Salvador and Gorka Guardiola Muzquiz. SealFs: Storage-based tamper-evident logging. *Comput. Secur.*, 2021.
- [42] John L Taylor and Tony Galica. A new code to protect victims in the uk from authorised push payments fraud. *Banking & Finance Law Review*, 2020.
- [43] The European Union Agency for Law Enforcement Cooperation. Take control of your digital life. don’t be a victim of cyber scams!, 2021. <https://www.europol.europa.eu/activities-services/public-awareness-and-prevention-guides/take-control-of-your-digital-life-don%E2%80%99t-be-victim-of-cyber-scams>.

- [44] The Federal Reserve. Fraud classifier. 2020. <https://fedpaymentsimprovement.org/wp-content/uploads/fraudclassifier-industry-adoption-roadmap.pdf>.
- [45] The Financial Ombudsman Service. Lending standards board review of the contingent reimbursement model code for authorised push payment scams-financial ombudsman service response. 2020. <https://www.financial-ombudsman.org.uk/files/289009/2020-10-02-LSB-CRM-Code-Review-Financial-Ombudsman-Service-Response.pdf>.
- [46] The International Criminal Police Organization. Investment fraud via dating apps, 2021. <https://www.interpol.int/en/News-and-Events/News/2021/Investment-fraud-via-dating-apps>.
- [47] UK Finance. 2021 half year fraud update, 2021. <https://www.ukfinance.org.uk/system/files/Half-year-fraud-update-2021-FINAL.pdf>.
- [48] UK Finance. The definitive overview of payment industry fraud, 2021. <https://www.ukfinance.org.uk/system/files/Fraud%20The%20Facts%202021-%20FINAL.pdf>.
- [49] Gavin Wood et al. Ethereum: A secure decentralised generalised transaction ledger. *Ethereum project yellow paper*, 2014.
- [50] Sophie Van Der Zee. Shifting the blame? investigation of user compliance with digital payment regulations. In *Cybercrime in Context*. Springer, 2021.

## A. Further Definition of APP fraud

An authorised push payment fraud has been defined by the “Financial Conduct Authority” (FCA) as: “*a transfer of funds by person A to person B, other than a transfer initiated by or through person B, where: (1) A intended to transfer the funds to a person other than B but was instead deceived into transferring the funds to B; or (2) A transferred funds to B for what they believed were legitimate purposes but which were, in fact, fraudulent*” [7]. Also, Payment Systems Regulator (PSR) defines APP fraud as follows: “*APP scams happen when a person or business is tricked into sending money to a fraudster posing as a genuine payee.*” [37].

## B. Commitment Scheme

A commitment scheme involves two parties, *sender* and *receiver*, and includes two phases: *commit* and *open*. In the commit phase, the sender commits to a message:  $x$  as  $\text{Com}(x, r) = \text{Com}_x$ , that involves a secret value:  $r \xleftarrow{\$} \{0, 1\}^\lambda$ . In the end of the commit phase, the commitment  $\text{Com}_x$  is sent to the receiver. In the open phase, the sender sends the opening  $\tilde{x} := (x, r)$  to the receiver who verifies its correctness:  $\text{Ver}(\text{Com}_x, \tilde{x}) \stackrel{?}{=} 1$  and accepts if the output is 1. A commitment scheme must satisfy two properties: (a) *hiding*: it is infeasible for an adversary (i.e., the receiver) to learn any information about the committed message  $x$ , until the commitment  $\text{Com}_x$  is opened, and (b) *binding*: it is infeasible for an adversary (i.e., the sender) to open a commitment  $\text{Com}_x$  to different values  $\tilde{x}' := (x', r')$  than that was used in the commit phase, i.e., infeasible to find  $\tilde{x}'$ , s.t.  $\text{Ver}(\text{Com}_x, \tilde{x}) = \text{Ver}(\text{Com}_x, \tilde{x}') = 1$ , where  $\tilde{x} \neq \tilde{x}'$ . There exist efficient non-interactive commitment schemes both in (a) the standard model, e.g., Pedersen scheme [38], and (b) the random oracle model using the well-known hash-based scheme such that committing is:  $H(x||r) = \text{Com}_x$  and  $\text{Ver}(\text{Com}_x, \tilde{x})$  requires checking:  $H(x||r) \stackrel{?}{=} \text{Com}_x$ , where  $H: \{0, 1\}^* \rightarrow \{0, 1\}^\lambda$  is a collision resistant hash function; i.e., the probability to find  $x$  and  $x'$  such that  $H(x) = H(x')$  is negligible in the security parameter  $\lambda$ .

## C. Bloom Filter

A Bloom filter [9] is a compact data structure for probabilistic efficient elements’ membership checking. A Bloom filter is an array of  $\bar{m}$  bits that are initially all set to zero. It represents  $\bar{n}$  elements. A Bloom filter comes along with  $\bar{k}$  independent hash functions. To insert an element, all the hash values of the element are computed and their corresponding bits in the filter are set to 1. To check an element’s membership, all its hash values are re-computed and checked whether all are set to one in the filter. If all the corresponding bits are one, then the element is probably in the filter; otherwise, it is not. In Bloom filters false positives are possible, i.e., it is possible that an element is not in the set, but the membership query shows that it is. According to [11], the upper bound of the false positive probability is:  $\bar{q} = \bar{p}^{\bar{k}} (1 + O(\frac{\bar{k}}{\bar{p}} \sqrt{\frac{\ln \bar{m} - \bar{k} \ln \bar{p}}{\bar{m}}}))$ , where  $\bar{p}$  is the probability that a particular bit in the filter is set to 1 and calculated as:  $\bar{p} = 1 - (1 - \frac{1}{\bar{m}})^{\bar{k}\bar{n}}$ . The efficiency of a Bloom filter depends on  $\bar{m}$  and  $\bar{k}$ . The lower bound of  $\bar{m}$  is  $\bar{n} \log_2 \bar{e} \cdot \log_2 \frac{1}{\bar{q}}$ , where  $\bar{e}$  is the base of natural logarithms, while the optimal number of hash functions is  $\log_2 \frac{1}{\bar{q}}$ , when  $\bar{m}$  is optimal. In this paper, we only use optimal  $\bar{k}$  and  $\bar{m}$ . In practice, we would like to have a predefined acceptable upper bound on false positive probability, e.g.,  $\bar{q} = 2^{-40}$ , and can adjust such a parameter. Given  $\bar{q}$  and  $\bar{n}$ , we can determine the rest of the parameters.

## D. Full Version of Definition of Payment with Dispute Resolution Scheme

In this section, we present a full formal definition of payment with dispute resolution. To allow this section to be self-contained, we repeat some content from Section 6. Below, we first provide the scheme’s syntax. Then, we formally define its correctness and security properties.

**Definition 7.** A payment with dispute resolution (pwwdr) involves six types of entities; namely, bank  $\mathcal{B}$ , customer  $\mathcal{C}$ , smart contract  $\mathcal{S}$ , certificate generator  $\mathcal{G}$ , set of auditors  $\mathcal{D} : \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$ , and dispute resolver  $\mathcal{DR}$ . It includes the following algorithms  $\text{pwwdr} := (\text{keyGen}, \text{bankInit}, \text{customerInit}, \text{genUpdateRequest}, \text{insertNewPayee}, \text{genWarning}, \text{genPaymentRequest}, \text{makePayment}, \text{genComplaint}, \text{verComplaint}, \text{resDispute})$ . Below, we define these algorithms.

- $\text{keyGen}(1^\lambda) \rightarrow (sk, pk)$ . It is a probabilistic algorithm run independently by  $\mathcal{G}$  and one of the auditors,  $\mathcal{D}_j$ . It takes as input a security parameter  $1^\lambda$ . It outputs a pair of secret keys  $sk := (sk_{\mathcal{G}}, sk_{\mathcal{D}})$  and public keys  $pk := (pk_{\mathcal{G}}, pk_{\mathcal{D}})$ , where  $sk_{\mathcal{D}}$  may contain multiple secret keys. The public key pair,  $pk$ , is sent to all participants.
- $\text{bankInit}(1^\lambda) \rightarrow (T, pp, l)$ . It is run by  $\mathcal{B}$ . It takes as input security parameter  $1^\lambda$ . It sets private parameters  $\tilde{\pi}_1$  and  $\tilde{\pi}_2$ . It generates an encoding-decoding token  $T$ , where  $T := (T_1, T_2)$ , each  $T_i$  contains a secret parameter  $\tilde{\pi}_i$  and its public witness  $g_i$ . Given a value and its witness anyone can check if they match. It also generates a set of (additional) public parameters,  $pp$ , one of which is  $e$  that is a threshold parameter. It also generates an empty list,  $l$ . It outputs  $T, pp$ , and  $l$ .  $\mathcal{B}$  sends  $(\tilde{\pi}_1, \tilde{\pi}_2)$  to  $\mathcal{C}$  and sends  $(g_1, g_2, pp, l)$  to  $\mathcal{S}$ .



- **customerInit**( $1^\lambda, T, pp$ )  $\rightarrow a$ . It is a deterministic algorithm run by  $C$ . It takes as input security parameter  $1^\lambda$ , token  $T$ , and set  $pp$  of public parameters. It checks the correctness of the elements in  $T$  and  $pp$ . If the checks pass, it outputs 1. Otherwise, it outputs 0.
- **genUpdateRequest**( $T, f, \mathbf{l}$ )  $\rightarrow \hat{m}_1^{(C)}$ . It is a deterministic algorithm run by  $C$ . It takes as input token  $T$ , new payee's detail  $f$ , and payees' list  $\mathbf{l}$ . It generates  $m_1^{(C)}$  which is an update request to the payees' list. It uses  $T_1 \in T$  and encoding algorithm  $\text{Encode}(T_1, \cdot)$  to encode  $m_1^{(C)}$  which results in  $\hat{m}_1^{(C)}$ . It outputs  $\hat{m}_1^{(C)}$ . Party  $C$  sends the output to  $S$ .
- **insertNewPayee**( $\hat{m}_1^{(C)}, \mathbf{l}$ )  $\rightarrow \hat{\mathbf{l}}$ . It is a deterministic algorithm run by  $S$ . It takes as input  $C$ 's encoded update request  $\hat{m}_1^{(C)}$ , and  $C$ 's payees' list  $\mathbf{l}$ . It inserts the new payee's detail into  $\mathbf{l}$  and outputs an updated list,  $\hat{\mathbf{l}}$ .
- **genWarning**( $T, \hat{\mathbf{l}}, aux$ )  $\rightarrow \hat{m}_1^{(B)}$ . It is run by  $B$ . It takes as input token  $T$ ,  $C$ 's encoded payees' list  $\hat{\mathbf{l}}$ , and auxiliary information:  $aux$ , e.g., a set of policies. Using  $T_1 \in T$ , it decodes and checks all elements of the list, e.g., whether they comply with the policies. If the check passes, it sets  $m_1^{(B)} = \text{"pass"}$ ; otherwise, it sets  $m_1^{(B)} = \text{warning}$ , where the warning is a string containing a warning detail along with the string "warning". It uses  $T_1$  and  $\text{Encode}(T_1, \cdot)$  to encode  $m_1^{(B)}$  which yields  $\hat{m}_1^{(B)}$ . It outputs  $\hat{m}_1^{(B)}$ . Party  $B$  sends  $\hat{m}_1^{(B)}$  to  $S$ .
- **genPaymentRequest**( $T, in_f, \hat{\mathbf{l}}, \hat{m}_1^{(B)}$ )  $\rightarrow \hat{m}_2^{(C)}$ . It is run by  $C$ . It takes as input token  $T$ , a payment detail  $in_f$ , encoded payees' list  $\hat{\mathbf{l}}$ , and encoded warning message,  $\hat{m}_1^{(B)}$ . Using  $T_1 \in T$ , it decodes  $\hat{\mathbf{l}}$  and  $\hat{m}_1^{(B)}$  yielding  $\mathbf{l}$  and  $m_1^{(B)}$  respectively. It checks the warning. It sets  $m_2^{(C)} = \phi$ , if it does not want to proceed. Otherwise, it sets  $m_2^{(C)}$  according to the content of  $in_f$  and  $\mathbf{l}$  (e.g., the amount of payment and payee's detail). It uses  $T_1$  and  $\text{Encode}(T_1, \cdot)$  to encode  $m_2^{(C)}$  resulting in  $\hat{m}_2^{(C)}$ . It outputs  $\hat{m}_2^{(C)}$ . Party  $C$  sends  $\hat{m}_2^{(C)}$  to  $S$ .
- **makePayment**( $T, \hat{m}_2^{(C)}$ )  $\rightarrow \hat{m}_2^{(B)}$ . It is a deterministic algorithm run by  $B$ . It takes as input token  $T$ , and encoded payment detail  $\hat{m}_2^{(C)}$ . Using  $T_1 \in T$ , it decodes  $\hat{m}_2^{(C)}$  and checks the result's validity, e.g., ensures it is well-formed or  $C$  has enough credit. If the check passes, it makes the payment and sets  $m_2^{(B)} = \text{"paid"}$ . Otherwise, it sets  $m_2^{(B)} = \phi$ . It uses  $T_1$  and  $\text{Encode}(T_1, \cdot)$  to encode  $m_2^{(B)}$  yielding  $\hat{m}_2^{(B)}$ . It outputs  $\hat{m}_2^{(B)}$ . Party  $B$  sends  $\hat{m}_2^{(B)}$  to  $S$ .
- **genComplaint**( $\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f$ )  $\rightarrow (\hat{z}, \hat{\pi})$ . It is run by  $C$ . It takes as input the encoded warning message  $\hat{m}_1^{(B)}$ , encoded payment message  $\hat{m}_2^{(B)}$ , token  $T$ , public key  $pk$ , and auxiliary information  $aux_f$ . It initially sets fresh strings  $(z_1, z_2, z_3)$  to null. Using  $T_1 \in T$ , it decodes  $\hat{m}_1^{(B)}$  and  $\hat{m}_2^{(B)}$  and checks the results' content. If it wants to complain that (i) "pass" message should have been a warning or (ii) no message was provided, it sets  $z_1$  to "challenge message". If its complaint is about the warning's effectiveness, it sets  $z_2$  to a combination of an evidence  $u \in aux_f$ , the evidence's certificate  $sig \in aux_f$ , the certificate's public parameter, and "challenge warning", where the certificate is obtained from  $\mathcal{G}$  via a query,  $Q$ . In certain cases, the certificate might be empty. If its complaint is about the payment, it sets  $z_3$  to "challenge payment". It uses  $T_1$  and  $\text{Encode}(T_1, \cdot)$  to encode  $z := (z_1, z_2, z_3)$  and uses  $pk_D$  and another encoding algorithm  $\text{Encode}(pk_D, \cdot)$  to encode  $\hat{\pi} := (\hat{\pi}_1, \hat{\pi}_2) \in T$ .

This results in  $\hat{z}$  and  $\hat{\pi}$  respectively. It outputs  $(\hat{z}, \hat{\pi})$ . Party  $C$  sends the pair to  $S$ .

- **verComplaint**( $\hat{z}, \hat{\pi}, g, \hat{\mathbf{m}}, \hat{\mathbf{l}}, j, sk_D, aux, pp$ )  $\rightarrow \hat{w}_j$ . It is run by every auditor  $\mathcal{D}_j$ . It takes as input  $C$ 's encoded complaint  $\hat{z}$ , encoded private parameters  $\hat{\pi}$ , the tokens' public parameters  $g := (g_1, g_2)$ , encoded messages  $\hat{\mathbf{m}} = [\hat{m}_1^{(C)}, \hat{m}_2^{(C)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ , encoded payees' list  $\hat{\mathbf{l}}$ , the auditor's index  $j$ , secret key set  $sk_D$ , auxiliary information  $aux$ , and set of public parameters  $pp$ . It uses a secret key in  $sk_D$  to decode  $\hat{\pi}$  that yields  $\hat{\pi} := (\hat{\pi}_1, \hat{\pi}_2)$ . It uses  $\hat{\pi}_1$  to decode  $\hat{z}, \hat{\mathbf{m}}$ , and  $\hat{\mathbf{l}}$  that results in  $z := (z_1, z_2, z_3), \mathbf{m}$ , and  $\mathbf{l}$  respectively. It checks if  $\hat{\pi}_i$  matches  $g_i$ . If the check fails, it aborts. It checks if  $m_1^{(C)}$  and  $m_2^{(C)}$  are non-empty; it aborts if the checks fail. It sets fresh parameters  $(w_{1,j}, w_{2,j}, w_{3,j}, w_{4,j})$  to 0. If "challenge message"  $\in z_1$ , given  $\mathbf{l}$ , it checks whether "pass" message (in  $m_1^{(B)}$ ) was given correctly or the missing message did not play any role in preventing the fraud. If either checks passes, it sets  $w_{1,j} = 0$ ; otherwise, it sets  $w_{1,j} = 1$ . If "challenge warning"  $\in z_2$ , it verifies the certificate in  $z_2$ . If it is invalid, it sets  $w_{3,j} = 0$ . If it is valid, it sets  $w_{3,j} = 1$ . It determines the effectiveness of the warning (in  $m_1^{(B)}$ ), by running an algorithm which determines that, i.e.,  $\text{checkWarning}(\cdot) \in aux$ . If it is effective, i.e.,  $\text{checkWarning}(m_1^{(B)}) = 1$ , it sets its verdict to 0, i.e.,  $w_{2,j} = 0$ ; otherwise, it sets  $w_{2,j} = 1$ . If "challenge payment"  $\in z_3$ , it checks if the payment was made (with the help of  $m_2^{(B)}$ ). If the check passes, it sets  $w_{4,j} = 1$ ; otherwise, it sets  $w_{4,j} = 0$ . It uses (another) secret key in  $sk_D$  and  $\hat{\pi}_2$  to encode  $\mathbf{w}_j = [w_{1,j}, w_{2,j}, w_{3,j}, w_{4,j}]$  yielding  $\hat{\mathbf{w}}_j = [\hat{w}_{1,j}, \hat{w}_{2,j}, \hat{w}_{3,j}, \hat{w}_{4,j}]$ . It outputs  $\hat{\mathbf{w}}_j$ . Party  $\mathcal{D}_j$  sends  $\hat{\mathbf{w}}_j$  to  $S$ .
- **resDispute**( $T_2, \hat{\mathbf{w}}, pp$ )  $\rightarrow \mathbf{v}$ . It is a deterministic algorithm run by  $DR$ . It takes as input token  $T_2$ , auditors' encoded verdicts  $\hat{\mathbf{w}} = [\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_n]$ , and public parameters set  $pp$ . It checks if the token's parameters match. If the check fails, it aborts. It uses  $\hat{\pi}_2 \in T_2$  to decode  $\hat{\mathbf{w}}$  and from the result it extracts final verdicts  $\mathbf{v} = [v_1, \dots, v_n]$ . The extraction procedure ensures each  $v_i$  is set to 1 only if at least  $e$  auditors' original verdicts (i.e.,  $w_{i,j}$ ) is 1, where  $e \in pp$ . It outputs  $\mathbf{v}$ . If  $v_4 = 1$  and (i) either  $v_1 = 1$  (ii) or  $v_2 = 1$  and  $v_3 = 1$ , then  $C$  is reimbursed.

Informally, pwdr has two properties; namely, *correctness* and *security*. Correctness requires that, in the absence of a fraudster, the payment journey is completed without the need for the honest customer to complain and the honest bank to reimburse the customer. Below, we formally state it.

**Definition 8** (Correctness). A pwdr scheme is correct if the key generation algorithm produces keys  $\text{keyGen}(1^\lambda) \rightarrow (sk, pk)$  such that for any payee's detail  $f$ , payment's detail  $in_f$ , and auxiliary information  $(aux, aux_f)$ , if  $\text{bankInit}(1^\lambda) \rightarrow (T, pp, \mathbf{l})$ ,  $\text{customerInit}(1^\lambda, T, pp) \rightarrow a$ ,  $\text{genUpdateRequest}(T, f, \mathbf{l}) \rightarrow \hat{m}_1^{(C)}$ ,  $\text{insertNewPayee}(\hat{m}_1^{(C)}, \mathbf{l}) \rightarrow \hat{\mathbf{l}}$ ,  $\text{genWarning}(T, \hat{\mathbf{l}}, aux) \rightarrow \hat{m}_1^{(B)}$ ,  $\text{genPaymentRequest}(T, in_f, \hat{\mathbf{l}}, \hat{m}_1^{(B)}) \rightarrow \hat{m}_2^{(C)}$ ,  $\text{makePayment}(T, \hat{m}_2^{(C)}) \rightarrow \hat{m}_2^{(B)}$ ,  $\text{genComplaint}(\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f) \rightarrow (\hat{z}, \hat{\pi})$ ,  $\forall j \in [n]$ :



$(\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{\mathbf{m}}, \hat{\mathbf{l}}, j, sk_D, aux, pp) \rightarrow \hat{\mathbf{w}}_j),$   
 $\text{resDispute}(T_2, \hat{\mathbf{w}}, pp) \rightarrow \mathbf{v}$ , then  
 $(z_1 = z_2 = z_3 = \phi) \wedge (\mathbf{v} = 0)$ , where  $g := (g_1, g_2) \in T$ ,  
 $\hat{\mathbf{m}} = [\hat{m}_1^{(c)}, \hat{m}_2^{(c)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ ,  $\hat{\mathbf{w}} = [\hat{w}_1, \dots, \hat{w}_n]$ , and  
 $z := (z_1, z_2, z_3)$  is the result of decoding  $\hat{z}$ .

A pwdr scheme is secure if it satisfies three main properties; namely, (a) security against a malicious victim, (b) security against a malicious bank, and (c) privacy. Below, we formally define them. Intuitively, security against a malicious victim requires that the victim of an APP fraud that is not qualified for reimbursement should not be reimbursed (despite it tries to be). More specifically, a corrupt victim cannot (a) make at least the threshold of the committee members,  $\mathcal{D}_j$ s, conclude that  $\mathcal{B}$  should have provided a warning, although  $\mathcal{B}$  has done so, or (b) make  $\mathcal{DR}$  conclude that the pass message was incorrectly given or a vital warning message was missing despite only less than the threshold of  $\mathcal{D}_j$ s believing so, or (c) persuade at least the threshold of  $\mathcal{D}_j$ s to conclude that the warning was ineffective although it was effective, or (d) make  $\mathcal{DR}$  believe that the warning message was ineffective although only less than the threshold of  $\mathcal{D}_j$ s believe that, or (e) convince  $\mathcal{D}_j$ s to accept an invalid certificate, or (f) make  $\mathcal{DR}$  believe that at least the threshold of  $\mathcal{D}_j$ s accepted the certificate although they did not, except for a negligible probability. Below, we formally state it.

**Definition 9** (Security against a malicious victim). A pwdr scheme is secure against a malicious victim, if for any security parameter  $\lambda$ , auxiliary information  $aux$ , and probabilistic polynomial-time adversary  $\mathcal{A}$ , there exists a negligible function  $\mu(\cdot)$ , such that for an experiment  $\text{Exp}_1^\lambda$ :

$\text{Exp}_1^\lambda(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow$  ( $sk, pk$ )
bankInit( $1^\lambda$ )  $\rightarrow$  ( $T, pp, \mathbf{l}$ )
 $\mathcal{A}(1^\lambda, T, pp, \mathbf{l}) \rightarrow \hat{m}_1^{(c)}$ 
insertNewPayee( $\hat{m}_1^{(c)}, \mathbf{l}$ )  $\rightarrow \hat{\mathbf{l}}$ 
genWarning( $T, \hat{\mathbf{l}}, aux$ )  $\rightarrow \hat{m}_1^{(B)}$ 
 $\mathcal{A}(T, \hat{\mathbf{l}}, \hat{m}_1^{(B)}) \rightarrow \hat{m}_2^{(c)}$ 
makePayment( $T, \hat{m}_2^{(c)}$ )  $\rightarrow \hat{m}_2^{(B)}$ 
 $\mathcal{A}(\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk) \rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n]$  :
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{\mathbf{m}}, \hat{\mathbf{l}}, j, sk_D, aux, pp) \rightarrow \hat{\mathbf{w}}_j$ )
resDispute( $T_2, \hat{\mathbf{w}}, pp$ )  $\rightarrow \mathbf{v} = [v_1, \dots, v_4]$ 

```

it holds that the following probability is  $\mu(\lambda)$ :

$$\Pr \left[ \begin{aligned} & \left( (m_1^{(B)} = \text{warning}) \wedge \left( \sum_{j=1}^n w_{1,j} \geq e \right) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{1,j} < e \right) \wedge (v_1 = 1) \right) \vee \\ & \left( (\text{checkWarning}(m_1^{(B)}) = 1) \wedge \right. \\ & \left. \left( \sum_{j=1}^n w_{2,j} \geq e \right) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{2,j} < e \right) \wedge (v_2 = 1) \right) \vee \\ & \left( u \notin Q \wedge \text{Sig.ver}(pk, u, sig) = 1 \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{3,j} < e \right) \wedge (v_3 = 1) \right) \end{aligned} \right] : \text{Exp}_1^\lambda(\text{input})$$

where  $g := (g_1, g_2) \in T$ ,  $\hat{\mathbf{m}} = [\hat{m}_1^{(c)}, \hat{m}_2^{(c)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ ,  $(w_{1,j}, \dots, w_{3,j})$  are the result of decoding  $(\hat{w}_{1,j}, \dots, \hat{w}_{3,j}) \in \hat{\mathbf{w}}_j \in \hat{\mathbf{w}}$ ,  $\text{checkWarning}(\cdot)$  determines a warning's effectiveness,  $\text{input} := (1^\lambda, aux)$ ,  $sk_D \in sk$ , and  $n$  is the total number of auditors. The probability is taken over the uniform choice of  $sk$ , randomness used in the blockchain's primitives (e.g., in signatures), randomness used during the encoding, and the randomness of  $\mathcal{A}$ .

Note, in the above experiment, we did not explicitly pass payment detail  $(f, in_f)$  and auxiliary information  $aux_f$  on to the adversary, to let the adversary pick them (on the client's behalf). Intuitively, security against a malicious bank requires that a malicious bank should not be able to disqualify an honest victim of an APP fraud from being reimbursed. Specifically, a corrupt bank cannot (a) make  $\mathcal{DR}$  conclude that the "pass" message was correctly given or an important warning message was not missing despite at least the threshold of  $\mathcal{D}_j$ s do not believe so, or (b) convince  $\mathcal{DR}$  that the warning message was effective although at least the threshold of  $\mathcal{D}_j$ s do not believe so, or (c) make  $\mathcal{DR}$  believe that less than the threshold of  $\mathcal{D}_j$ s did not accept the certificate although at least the threshold of them did that, or (d) make  $\mathcal{DR}$  believe that no payment was made, although at least the threshold of  $\mathcal{D}_j$ s believe the opposite, except for a negligible probability. Below, we formally state it.

**Definition 10** (Security against a malicious bank). A pwdr scheme is secure against a malicious bank, if for any security parameter  $\lambda$ , auxiliary information  $aux$ , and probabilistic polynomial-time adversary  $\mathcal{A}$ , there exists a negligible function  $\mu(\cdot)$ , such that for an experiment  $\text{Exp}_2^\lambda$ :

$\text{Exp}_2^\lambda(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow$  ( $sk, pk$ )
 $\mathcal{A}(1^\lambda) \rightarrow (T, pp, \mathbf{l}, f, in_f, aux_f)$ 
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
genUpdateRequest( $T, f, \mathbf{l}$ )  $\rightarrow \hat{m}_1^{(c)}$ 
insertNewPayee( $\hat{m}_1^{(c)}, \mathbf{l}$ )  $\rightarrow \hat{\mathbf{l}}$ 
 $\mathcal{A}(T, \hat{\mathbf{l}}, aux) \rightarrow \hat{m}_1^{(B)}$ 
genPaymentRequest( $T, in_f, \hat{\mathbf{l}}, \hat{m}_1^{(B)}$ )  $\rightarrow \hat{m}_2^{(c)}$ 
 $\mathcal{A}(T, \hat{m}_2^{(c)}) \rightarrow \hat{m}_2^{(B)}$ 
genComplaint( $\hat{m}_1^{(B)}, \hat{m}_2^{(B)}, T, pk, aux_f$ )  $\rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n]$  :
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{\mathbf{m}}, \hat{\mathbf{l}}, j, sk_D, aux, pp) \rightarrow \hat{\mathbf{w}}_j$ )
resDispute( $T_2, \hat{\mathbf{w}}, pp$ )  $\rightarrow \mathbf{v} = [v_1, \dots, v_4]$ 

```

it holds that:

$$\Pr \left[ \begin{aligned} & \left( \left( \sum_{j=1}^n w_{1,j} \geq e \right) \wedge (v_1 = 0) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{2,j} \geq e \right) \wedge (v_2 = 0) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{3,j} \geq e \right) \wedge (v_3 = 0) \right) \vee \\ & \left( \left( \sum_{j=1}^n w_{4,j} \geq e \right) \wedge (v_4 = 0) \right) \end{aligned} \right] : \text{Exp}_2^\lambda(\text{input})$$

where  $g := (g_1, g_2) \in T$ ,  $\hat{\mathbf{m}} = [\hat{m}_1^{(c)}, \hat{m}_2^{(c)}, \hat{m}_1^{(B)}, \hat{m}_2^{(B)}]$ ,  $(w_{1,j}, \dots, w_{4,j})$  are the result of decoding  $(\hat{w}_{1,j}, \dots, \hat{w}_{4,j}) \in \hat{\mathbf{w}}_j \in \hat{\mathbf{w}}$ ,  $\text{input} := (1^\lambda, aux)$ ,  $sk_D \in sk$ ,  $n$  is the total number of auditors. The probability is taken over the uniform choice of  $sk$ , randomness used in the blockchain's

primitives, randomness used during the encoding, and the randomness of  $\mathcal{A}$ .

Now we move on to privacy. Informally, a pwdr is privacy-preserving if it protects the privacy of (1) the customers', bank's, and auditors' messages (except public parameters) from non-participants of the scheme, including other customers, and (2) each auditor's verdict from  $\mathcal{DR}$  which sees the final verdict. Below, we formally state it.

**Definition 11 (Privacy).** A pwdr preserves privacy if the following two properties are satisfied.

- 1) For any probabilistic polynomial-time adversary  $\mathcal{A}_1$ , security parameter  $\lambda$ , and auxiliary information  $aux$ , there exists a negligible function  $\mu(\cdot)$ , such that for any experiment  $\text{Exp}_3^{\mathcal{A}_1}$ :

$\text{Exp}_3^{\mathcal{A}_1}(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow$  ( $sk, pk$ )
bankInit( $1^\lambda$ )  $\rightarrow$  ( $T, pp, l$ )
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
 $\mathcal{A}_1(1^\lambda, pk, a, pp, g, l) \rightarrow ((f_0, f_1), (in_{f_0}, in_{f_1}),$ 
 $(aux_{f_0}, aux_{f_1}))$ 
 $\gamma \xleftarrow{\$} \{0, 1\}$ 
genUpdateRequest( $T, f_\gamma, l$ )  $\rightarrow \hat{m}_1^{(c)}$ 
insertNewPayee( $\hat{m}_1^{(c)}, l$ )  $\rightarrow \hat{l}$ 
genWarning( $T, \hat{l}, aux$ )  $\rightarrow \hat{m}_1^{(b)}$ 
genPaymentRequest( $T, in_{f_\gamma}, \hat{l}, \hat{m}_1^{(b)}$ )  $\rightarrow \hat{m}_2^{(c)}$ 
makePayment( $T, \hat{m}_2^{(c)}$ )  $\rightarrow \hat{m}_2^{(b)}$ 
genComplaint( $\hat{m}_1^{(b)}, \hat{m}_2^{(b)}, T, pk, aux_{f_\gamma}$ )  $\rightarrow (\hat{z}, \hat{\pi})$ 
 $\forall j, j \in [n]:$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v$ 

```

it holds that the following probability is smaller than equal to  $\frac{1}{2} + \mu(\lambda)$ .

$$\Pr[\mathcal{A}_1(g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w}) \rightarrow \gamma : \text{Exp}_3^{\mathcal{A}_1}(\text{input})]$$

- 2) For any probabilistic polynomial-time adversaries  $\mathcal{A}_2$  and  $\mathcal{A}_3$ , security parameter  $\lambda$ , and auxiliary information  $aux$ , there exists a negligible function  $\mu(\cdot)$ , such that for any experiment  $\text{Exp}_4^{\mathcal{A}_2}$ :

$\text{Exp}_4^{\mathcal{A}_2}(1^\lambda, aux)$

```

keyGen( $1^\lambda$ )  $\rightarrow$  ( $sk, pk$ )
bankInit( $1^\lambda$ )  $\rightarrow$  ( $T, pp, l$ )
customerInit( $1^\lambda, T, pp$ )  $\rightarrow a$ 
 $\mathcal{A}_2(1^\lambda, pk, a, pp, l) \rightarrow (f, in_f, aux_f)$ 
genUpdateRequest( $T, f, l$ )  $\rightarrow \hat{m}_1^{(c)}$ 
insertNewPayee( $\hat{m}_1^{(c)}, l$ )  $\rightarrow \hat{l}$ 
 $\mathcal{A}_2(T, \hat{l}, aux) \rightarrow m_1^{(b)}$ 
Encode( $T_1, m_1^{(b)}$ )  $\rightarrow \hat{m}_1^{(b)}$ 
genPaymentRequest( $T, in_f, \hat{l}, \hat{m}_1^{(b)}$ )  $\rightarrow \hat{m}_2^{(c)}$ 
 $\mathcal{A}_2(T, pk, aux_f, \hat{m}_1^{(b)}, \hat{m}_2^{(c)}) \rightarrow (m_2^{(b)}, z, \hat{\pi})$ 
Encode( $T_1, m_2^{(b)}$ )  $\rightarrow \hat{m}_2^{(b)}$ 
Encode( $T_1, z$ )  $\rightarrow \hat{z}$ 
Encode( $pk_D, \hat{\pi}$ )  $\rightarrow \hat{\pi}$ 
 $\forall j, j \in [n]:$ 
  ( $\text{verComplaint}(\hat{z}, \hat{\pi}, g, \hat{m}, \hat{l}, j, sk_D, aux, pp) \rightarrow \hat{w}_j$ )
resDispute( $T_2, \hat{w}, pp$ )  $\rightarrow v$ 

```

it holds that:

$$\Pr[\mathcal{A}_3(T_2, pk, pp, g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w}, v) \rightarrow w_j : \text{Exp}_4^{\mathcal{A}_2}(\text{input})] \leq \Pr' + \mu(\lambda),$$

where  $g := (g_1, g_2) \in T$ ,  $\hat{m} = [\hat{m}_1^{(c)}, \hat{m}_2^{(c)}, \hat{m}_1^{(b)}, \hat{m}_2^{(b)}]$ ,  $\hat{w} = [\hat{w}_1, \dots, \hat{w}_n]$ ,  $\text{input} := (1^\lambda, aux)$ ,  $T_1 \in T$ ,  $pk_D \in pk$ ,  $sk_D \in sk$ ,  $n$  is the total number of auditors. Moreover,  $\Pr'$  is defined as follows. Let auditor  $\mathcal{D}_i$  output 0 and 1 with probabilities  $\Pr_{i,0}$  and  $\Pr_{i,1}$  respectively. Then,  $\Pr'$  is defined as  $\text{Max}(\Pr_{1,0}, \Pr_{1,1}, \dots, \Pr_{n,0}, \Pr_{n,1})$ . In the above privacy definition, the probability is taken over the uniform choice of  $sk$ , the probability that each  $\mathcal{D}_j$  outputs 0 or 1, the randomness used in the blockchain's primitives, the randomness used during the encoding, and the randomness of  $\mathcal{A}_1$  in  $\text{Exp}_3^{\mathcal{A}_1}$  and  $\mathcal{A}_2$  in  $\text{Exp}_4^{\mathcal{A}_2}$ .

**Definition 12 (Security).** A pwdr is secure if it meets security against a malicious victim, security against a malicious bank, and preserves privacy with respect to definitions 9, 10, and 11 respectively.

## E. Further Discussion on Privacy Definition

In this section, we discuss the intuition behind the privacy definition, i.e., Definition 5. Below, we explain each case.

### E.1. Case 1

Overall, the idea behind the design of experiment  $\text{Exp}_3^{\mathcal{A}_1}(\cdot)$  is similar to that of the semantic security of an encryption scheme [26]. In this experiment, an adversary picks plaintext inputs:  $(f_0, f_1), (in_{f_0}, in_{f_1}), (aux_{f_0}, aux_{f_1})$ . Then, the experiment (or oracle) flips a coin to pick a random index,  $\gamma \xleftarrow{\$} \{0, 1\}$ , and uses only one element from each pair, i.e.,  $(f_\gamma, in_{f_\gamma}, aux_{f_\gamma})$ , and honestly executes all algorithms of PwDR (given the chosen inputs). In this case, the privacy property states that given the public messages  $(g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w})$ , adversary  $\mathcal{A}_1$  should not be able to tell which index the oracle picked (i.e., find the value of  $\gamma$ ), with a probability significantly better than  $\frac{1}{2}$ ; more formally its probability of success should be at most  $\frac{1}{2} + \mu(\lambda)$ .

### E.2. Case 2

In this case, we allow adversary  $\mathcal{A}_2$  to arbitrarily pick the plaintext inputs (on behalf of the client and bank) and feed them to the PwDR algorithms. In this experiment,  $\mathcal{A}_2$  can craft its input in a way that an arbiter's verdict is in favour of the client or the bank, i.e., 1 or 0. For instance, the adversary can generate an invalid claim on behalf of the client; in this case, each arbiter's verdict will be against the client, i.e.,  $w_j = 0$ .

Therefore, to capture  $\mathcal{A}_2$ 's arbitrary behaviour in picking valid or invalid inputs, we let auditor  $\mathcal{D}_i$  output 0 and

1 with probabilities  $Pr_{i,0}$  and  $Pr_{i,1}$  respectively. Then, we define  $Pr'$  as  $Max(Pr_{1,0}, Pr_{1,1}, \dots, Pr_{n,0}, Pr_{n,1})$ .

In this case, the privacy property states that given the public messages (and the messages that a dispute resolver receives) another adversary  $\mathcal{A}_3$  (which does not interact with  $\mathcal{A}_2$ ) should not be able to tell the value of each auditor's verdict  $w_j$  with a probability significantly better than  $Pr'$ ; more formally its probability of success should be at most  $Pr' + \mu(\lambda)$ .

Note that, in  $\text{Exp}_4^{\mathcal{A}_2}(\cdot)$ , we explicitly invoked  $\text{Encode}(\cdot)$  and  $\text{Encode}(\cdot)$  (defined in the PwDR syntax in Section 6) to ensure that the messages that  $\mathcal{A}_2$  generates are correctly encoded; otherwise, if we allowed  $\mathcal{A}_2$  to encode the messages, then the privacy notion would not make sense as it could simply avoid doing so.

## F. Further Discussion on Variant 1 Encoding-Decoding Protocol

In this section, we first formally state our observation on which Variant 1 encoding-decoding protocol relies and then prove it. After that, we explain why this variant meets the three properties we laid out in Section 7.4, i.e., it should (1) generate unlinkable verdicts, (2) not require auditors to interact with each other for each customer, and (3) be efficient.

**Theorem 2.** *Let set  $S = \{s_1, \dots, s_m\}$  be the union of two disjoint sets  $S'$  and  $S''$ , where  $S'$  contains non-zero random values picked uniformly from a finite field  $\mathbb{F}_p$ ,  $S''$  contains zeros,  $|S'| \geq c' = 1$ ,  $|S''| \geq c'' = 0$ , and pair  $(c', c'')$  is public information. Then,  $r = \bigoplus_{i=1}^m s_i$  reveals nothing beyond the public information.*

*Proof.* Let  $s_1$  and  $s$ , be two random values picked uniformly at random from  $\mathbb{F}_p$ . Let  $\bar{s} = s_1 \oplus \underbrace{0 \oplus \dots \oplus 0}_{|S''|}$ .

Since  $\bar{s} = s_1$ , two values  $\bar{s}$  and  $s$  have identical distribution. Thus,  $\bar{s}$  reveals nothing in this case. Next, let  $\tilde{s} = \underbrace{s_1 \oplus s_2 \oplus \dots \oplus s_j}_{|S'|}$ , where  $s_i \in S'$ . Since each  $s_i$

is a uniformly random value, the XOR of them is a uniformly random value too. That means values  $\tilde{s}$  and  $s$  have identical distributions. Thus,  $\tilde{s}$  reveals nothing in this case as well. Also, it is not hard to see that the combination of the above two cases reveals nothing too, i.e.,  $\bar{s} \oplus \tilde{s}$  and  $s$  have identical distribution.  $\square$

The primary reason this variant meets property 1 is that each masked verdict reveals nothing about the verdict (and its representation) and given the final verdict,  $\mathcal{I}$  cannot distinguish between the case where there is exactly one auditor that voted 1 and the case where multiple auditors voted 1, as in both cases  $\mathcal{I}$  extracts only a single random value, which reveals nothing about the number of auditors which voted 0 or 1 (due to Theorem 2). Note, the protocols' correctness holds with an overwhelming probability, i.e.,  $1 - \frac{1}{2^\lambda}$ . Specifically, if two auditors represent their verdict by an identical random value, then when they are XORed they would cancel out each other which can affect the result's correctness. The same holds if the XOR of multiple verdicts' representations results in a value that can cancel out another verdict's representation.

Nevertheless, the probability that such an event occurs is negligible in the security parameter  $|p| = \lambda$ , i.e., the probability is at most  $\frac{1}{2^\lambda}$ . It is evident that this variant meets property 2 as the auditors interact with each other *only once* (to agree on a key) for all customers. It also meets property 3 as it involves pseudorandom function invocations and XOR operations which are highly efficient operations.

## G. Generic Verdict Encoding-Decoding Protocol

In this section, we present the generic verdict encoding-decoding protocols (i.e., GPVE and GFVD) in detail, in Figures 6 and 7. They let a semi-honest third party  $\mathcal{I}$  find out if at least  $e$  auditors voted 1, where  $e$  can be any integer in the range  $[1, n]$ . As Figure 6 indicates (and as we discussed in Section 7.4) after  $\mathcal{D}_n$  generates all combinations of verdict 1's representations, it inserts the combinations into a Bloom filter, to preserve the representations' privacy from  $\mathcal{I}$ . Note, instead of inserting the combinations into a Bloom filter, we could *hash* the combinations and give the hash values to  $\mathcal{I}$ . However, using a Bloom filter lets us save considerable communication costs. Let us see a concrete example. Let  $n = 10$ ,  $e = 6$ , and the hash function be SHA-256. If the latter (hash-based) approach is used,  $\mathcal{D}_n$  needs to send  $|W| \times 256 = 386 \times 256 = 98,816$  bits to  $\mathcal{I}$ , whereas if the former (Bloom filter-based) approach is used, then it only needs to send  $|\text{BF}| = 22,276$  bits to  $\mathcal{I}$ . Thus, by using a Bloom filter, it can save communication costs by at least a factor of 4.

## H. Variant 2 Encoding-Decoding Protocol's Main Theorem and Proof

In this section, we first formally state our main observation on which Variant 2 encoding-decoding protocol relies. After that, we prove it.

**Theorem 3.** *Let  $S = \{s_1, \dots, s_m\}$  be a set of random values picked uniformly from finite field  $\mathbb{F}_p$ , where the cardinality of  $S$  is public information. Let BF be a Bloom filter encoding all elements of  $S$ . Then, BF reveals nothing about any element of  $S$ , beyond the public information, except with a negligible probability in the security parameter  $\lambda$ , i.e., with a probability at most  $\frac{|S|}{2^\lambda}$ .*

*Proof.* First, we consider the simplest case where only a single element of  $S$  is encoded in BF. In this case, due to the pre-image resistance of the Bloom filter's hash functions and the fact that the set's element was picked uniformly at random from  $\mathbb{F}_p$ , the probability that BF reveals anything about the original element is at most  $\frac{1}{2^\lambda}$ . Now, we move on to the case where all elements of  $S$  are encoded in BF. In this case, the probability that BF reveals anything about at least an element of the set is  $\frac{|S|}{2^\lambda}$ , due to the pre-image resistance of the hash functions, the fact that all elements were selected uniformly at random from the finite field, and the union bound. Nevertheless, when a BF's size is set appropriately to avoid false-positive without wasting storage, this reveals the number of elements encoded in it, which is public

GPVE( $\bar{k}_0, \text{ID}, w_j, o, e, n, j$ )  $\rightarrow (\bar{w}_j, \text{BF})$

- *Input.*  $\bar{k}_0$ : a key of pseudorandom function PRF( $\cdot$ ), ID: a unique identifier,  $w_j$ : a verdict,  $o$ : a counter,  $e$ : a threshold,  $n$ : the total number of auditors, and  $j$ : an auditor's index.

- *Output.*  $\bar{w}_j$ : an encoded verdict.

Auditor  $\mathcal{D}_j$  takes the following steps.

- 1) computes a pseudorandom value, as follows.
  - if  $j < n$  :  $r_j = \text{PRF}(\bar{k}_0, 1 || o || j || \text{ID})$ .

- if  $j = n$  :  $r_j = \bigoplus_{i=1}^{n-1} r_i$ .

Note, the above second step is taken only by  $\mathcal{D}_n$ .

- 2) sets a fresh parameter,  $w'_j$ , that represents a verdict, as below.

$$w'_j = \begin{cases} 0, & \text{if } w_j = 0 \\ \alpha_j = \text{PRF}(\bar{k}_0, 2 || o || j || \text{ID}), & \text{if } w_j = 1 \end{cases}$$

- 3) masks  $w'_j$  as follows.  $\bar{w}_j = w'_j \oplus r_j$ .
- 4) if  $j = n$ , computes a Bloom filter that encodes the combinations of verdict representations (i.e.,  $w'_j$ ) for verdict 1. In particular, it takes the following steps.
  - for every integer  $i$  in the range  $[e, n]$ , computes the combinations (without repetition) of  $i$  elements from set  $\{\alpha_1, \dots, \alpha_n\}$ . In the case where multiple elements are taken at a time (i.e.,  $i > 1$ ), the elements are XORed with each other. Let  $W = \{(\alpha_1 \oplus \dots \oplus \alpha_e), (\alpha_2 \oplus \dots \oplus \alpha_{e+1}), \dots, (\alpha_1 \oplus \dots \oplus \alpha_n)\}$  be the result.
  - constructs an empty Bloom filter. Then, it inserts all elements of  $W$  into this Bloom filter. Let BF be the Bloom filter encoding  $W$ 's elements.
- 5) outputs  $(\bar{w}_j, \text{BF})$ .

Figure 6: Generic Private Verdict Encoding (GPVE) Protocol. In the figure,  $\mathcal{D}_n$  can generate other auditors'  $r_i$  and  $\alpha_j$  values, given  $\bar{k}_0$ . Note also that ID is a unique identifier (e.g., wallet address) of the party for whom a verdict is provided (e.g., a client), and  $o$  is a counter that determines how many times a verdict for the same ID holder has been generated in the past. ID and  $o$  are used to ensure each  $r_j$  will be different for each invocation of GPVE although the same key  $\bar{k}_0$  is used.

GFVD( $n, \bar{w}, \text{BF}$ )  $\rightarrow v$

- *Input.*  $n$ : the total number of auditors, and  $\bar{w} = [\bar{w}_1, \dots, \bar{w}_n]$ : a vector of all auditors' encoded verdicts.
- *Output.*  $v$ : final verdict.

A third-party  $\mathcal{I}$  takes the following steps.

- 1) combines all auditors' encoded verdicts,  $\bar{w}_j \in \bar{w}$ , as follows.  $c = \bigoplus_{j=1}^n \bar{w}_j$
- 2) checks if  $c$  is in the Bloom filter, BF.
- 3) sets the final verdict  $v$  depending on the content of  $c$ . Specifically,

$$v = \begin{cases} 0, & \text{if } c = 0 \text{ or } c \notin \text{BF} \\ 1, & \text{if } c \in \text{BF} \end{cases}$$

- 4) outputs  $v$ .

Figure 7: Generic Final Verdict Decoding (GFVD) Protocol.

information. Thus, the only information BF reveals is the public one.  $\square$

## I. Further Discussion on the Verdict Encoding-decoding Protocols

Recall that each variant of our verdict encoding-decoding protocol is a voting mechanism. It lets a third party,  $\mathcal{I}$ , find out if a threshold of the auditors voted 1, while (i) generating unlinkable verdicts, (ii) not requiring auditors to interact with each other for each customer, (iii) hiding the number of 0 or 1 verdicts from  $\mathcal{I}$ , and (iv) being efficient. Therefore, it is natural to ask:

*Is there any e-voting protocol, in the literature, that can simultaneously satisfy all the above requirements?*

The short answer is no. Recently, a provably secure e-voting protocol that can hide the number of 1 and 0 votes has been proposed by Kusters *et al.* [28]. Although this scheme can satisfy the above security requirements, it imposes a high computation cost, as it involves computationally expensive primitives such as zero-knowledge proofs, threshold public-key encryption scheme, and generic multi-party computation. In contrast, our verdict encoding-decoding protocols rely on much more lightweight operations such as XOR and hash function evaluations. We also highlight that our verdict encoding-decoding protocols are in a different setting than the one in which most of the e-voting protocols are. Because the former protocols are in the setting where there exists a small number of auditors (or voters) which are trusted and can interact with each other once; whereas, the latter (e-voting) protocols are in a more generic setting where there is a large number of voters, some of which might be malicious, and they are not required to interact with each other.

Note that each variant of our verdict encoding-decoding protocol requires every auditor to provide an encoded vote in order for  $\mathcal{I}$  to extract the final verdict. To let each variant terminate and  $\mathcal{I}$  find out the final verdict in the case where a set of auditors do not provide their vote, we can integrate the following idea into each variant. We define a manager auditor, say  $\mathcal{D}_n$ , which is always responsive and keeps track of missing votes. After the voting time elapses and  $\mathcal{D}_n$  realises a certain number of auditors did not provide their encoded vote, it provides 0 votes on their behalf and masks them using the auditors' masking values.

## J. Security Analysis of the PwDR Protocol

In this section, we prove the security theorem of the PwDR protocol, i.e., Theorem 1. To prove this theorem, we show that the PwDR satisfies all security properties defined in Section 6. We first prove that it meets security against a malicious victim.

**Lemma 1.** *If the digital signature is existentially unforgeable under chosen message attacks, and the SAP and blockchain are secure, then the PwDR scheme is secure against a malicious victim, with regard to Definition 9.*

*Proof.* First, we focus on event I :  $\left( (m_1^{(B)} = \text{warning}) \wedge \left( \sum_{j=1}^n w_{1,j} \geq e \right) \right)$  which considers the case where  $\mathcal{B}$  has provided a warning message but  $\mathcal{C}$  manages to convince at least the threshold of the auditors to set their verdicts to 1, that ultimately results in  $\sum_{j=1}^n w_{1,j} \geq e$ . We argue that the adversary's success probability in this event is negligible in the security parameter. Specifically, due to the security of SAP (i.e., the binding property of the SAP's commitment),  $\mathcal{C}$  cannot convince an auditor to accept a different decryption key, e.g.,  $k' \in \tilde{\pi}'$ , that will be used to decrypt  $\mathcal{B}$ 's encrypted message  $\hat{m}_1^{(B)}$ , other than what was agreed between  $\mathcal{C}$  and  $\mathcal{B}$  in the initiation phase, i.e.,  $\bar{k}_1 \in \tilde{\pi}_1$ . To be more precise, it cannot persuade an auditor to accept a statement  $\tilde{\pi}'$ , where  $\tilde{\pi}' \neq \tilde{\pi}_1$ , except with a negligible probability,  $\mu(\lambda)$ . This ensures that honest  $\mathcal{B}$ 's original message (and accordingly the warning) is accessed by every auditor with a high probability. Next, we consider event II :  $\left( \left( \sum_{j=1}^n w_{1,j} < e \right) \wedge (v_1 = 1) \right)$  that captures the case where only less than the threshold of the auditors approved that the pass message was given incorrectly or the missing message could prevent the APP fraud, but the final verdict that  $\mathcal{DR}$  extracts implies that at least the threshold of the auditors approved that. We argue that the probability that this event occurs is negligible in the security parameter too. Specifically, due to the binding property of the SAP,  $\mathcal{C}$  cannot persuade (a) an auditor to accept a different encryption key and (b)  $\mathcal{DR}$  to accept a different decryption key other than what was agreed between  $\mathcal{C}$  and  $\mathcal{B}$  in the initiation phase. More precisely, it cannot persuade them to accept a statement  $\tilde{\pi}'$ , where  $\tilde{\pi}' \neq \tilde{\pi}_2$ , except with a negligible probability,  $\mu(\lambda)$ .

Now, we move on to event III :  $\left( (\text{checkWarning}(m_1^{(B)}) = 1) \wedge \left( \sum_{j=1}^n w_{2,j} \geq e \right) \right)$ . It captures the case where  $\mathcal{B}$  has provided an effective warning message but  $\mathcal{C}$  manages to make at least the threshold of the auditors set their verdicts to 1, that ultimately results in  $\sum_{j=1}^n w_{2,j} \geq e$ . The same argument provided to event I is applicable to this event too. Briefly, due to the security of the SAP,  $\mathcal{C}$  cannot persuade an auditor to accept a different decryption key other than what was agreed between  $\mathcal{C}$  and  $\mathcal{B}$  in the initiation phase. Therefore, all auditors will receive the original message of  $\mathcal{B}$ , including the effective warning message, except a negligible probability,  $\mu(\lambda)$ . Now, we consider event IV :  $\left( \left( \sum_{j=1}^n w_{2,j} < e \right) \wedge (v_2 = 1) \right)$  which captures the case where at least the threshold of the auditors approved that the warning message was effective but the final verdict that  $\mathcal{DR}$  extracts implies that they approved the opposite. The security argument of event II applies to this event as well. In short, due to the security of the SAP,  $\mathcal{C}$  cannot persuade an auditor to accept a different encryption key, and cannot convince  $\mathcal{DR}$  to accept a different decryption key other than what was initially agreed between  $\mathcal{C}$  and  $\mathcal{B}$ , except a negligible probability,  $\mu(\lambda)$ . Now, we analyse event V :  $\left( u \notin Q \wedge \text{Sig.ver}(pk, u, sig) = 1 \right)$ . This event captures the case where the malicious victim comes up with a valid signature/certificate on a message

that has never been queried to the signing oracle. But, due to the existential unforgeability of the digital signature scheme, the probability that such an event occurs is negligible,  $\mu(\lambda)$ . Next, we focus on event VI :  $\left( \left( \sum_{j=1}^n w_{3,j} < e \right) \wedge (v_3 = 1) \right)$  that considers the case where less than the threshold of the auditors indicated that the signature (in  $\mathcal{C}$ 's complaint) is valid, but the final verdict that  $\mathcal{DR}$  extracts implies that at least the threshold of the auditors approved the signature. This means the adversary has managed to switch the verdicts of those auditors who voted 0 to 1. However, the probability that this even occurs is negligible as well. Because, due to the SAP's binding property,  $\mathcal{C}$  cannot convince an auditor and  $\mathcal{DR}$  to accept different encryption and decryption keys other than what was initially agreed between  $\mathcal{C}$  and  $\mathcal{B}$ , except a negligible probability,  $\mu(\lambda)$ . Therefore, with only a negligible probability the adversary can switch a verdict for 0 to the verdict for 1.

Moreover, a malicious  $\mathcal{C}$  cannot frame an honest  $\mathcal{B}$  for providing an invalid message by manipulating the smart contract's content, e.g., by replacing an effective warning with an ineffective one in  $\mathcal{S}$ , or excluding a warning from  $\mathcal{S}$ . To do that, it has to either forge the honest party's signature, so it can send an invalid message on its behalf, or fork the blockchain so the chain comprising a valid message is discarded. In the former case, the adversary's probability of success is negligible as long as the signature is secure. The adversary has the same success probability in the latter case because it has to generate a long enough chain that excludes the valid message which has a negligible success probability, under the assumption that the hash power of the adversary is lower than those of honest miners and due to the blockchain's liveness property an honestly generated transaction will eventually appear on an honest miner's chain [23].  $\square$

Now, we first present a lemma formally stating that the PwDR protocol is secure against a malicious bank and then prove this lemma.

**Lemma 2.** *If the SAP and blockchain are secure, and the correctness of verdict encoding-decoding protocols (i.e., PVE and FVD) holds, then the PwDR protocol is secure against a malicious bank, with regard to Definition 10.*

*Proof.* We first focus on event I :  $\left( \left( \sum_{j=1}^n w_{1,j} \geq e \right) \wedge (v_1 = 0) \right)$  which captures the case where  $\mathcal{DR}$  is convinced that the pass message was correctly given or an important warning message was not missing, despite at least the threshold of the auditors do not believe so. We argue that the probability that this event takes place is negligible in the security parameter. Because, due to the SAP's binding property,  $\mathcal{B}$  cannot persuade  $\mathcal{DR}$  to accept a different decryption key, e.g.,  $k' \in \tilde{\pi}'$ , other than what was agreed between  $\mathcal{C}$  and  $\mathcal{B}$  in the initiation phase, i.e.,  $\bar{k}_2 \in \tilde{\pi}_2$ , except with a negligible probability. Specifically, it cannot persuade  $\mathcal{DR}$  to accept a statement  $\tilde{\pi}'$ , where  $\tilde{\pi}' \neq \tilde{\pi}_2$  except with probability  $\mu(\lambda)$ . Also, as discussed in Section 7.4, due to the correctness of the verdict encoding-decoding protocols, i.e., PVE and FVD, the probability that multiple representations of verdict 1 cancel out each other is negligible too, i.e., it is at

most  $\frac{1}{2^\lambda}$ . Thus, event I occurs only with a negligible probability,  $\mu(\lambda)$ . To assert that events II :  $\left(\left(\sum_{j=1}^n w_{2,j} \geq e\right) \wedge (v_2 = 0)\right)$ , III :  $\left(\left(\sum_{j=1}^n w_{3,j} \geq e\right) \wedge (v_3 = 0)\right)$ , and IV :  $\left(\left(\sum_{j=1}^n w_{4,j} \geq e\right) \wedge (v_4 = 0)\right)$  occur only with a negligible probability, we can directly use the above argument provided for event I. To avoid repetition, we do not restate them in this proof. Moreover, a malicious  $\mathcal{B}$  cannot frame an honest  $\mathcal{C}$  for providing an invalid message by manipulating the smart contract's content, e.g., by replacing its valid signature with an invalid one or sending a message on its behalf, due to the security of the blockchain.  $\square$

Next, we prove the PwDR protocol's privacy. As before, we first formally state the related lemma and then prove it.

**Lemma 3.** *If the encryption schemes are semantically secure, and the SAP and encoding-decoding schemes (i.e., PVE and FVD) are secure, then the PwDR protocol is privacy-preserving with regard to Definition 11.*

*Proof.* We first focus on property 1, i.e., the privacy of the parties' messages from the public. Due to the privacy-preserving property of the SAP, that relies on the hiding property of the commitment scheme, given the public commitments,  $g := (g_1, g_2)$ , the adversary learns no information about the committed values,  $(\bar{k}_1, \bar{k}_2)$ , except with a negligible probability,  $\mu(\lambda)$ . Thus, it cannot find the encryption-decryption keys used to generate ciphertext  $\hat{m}, \hat{l}, \hat{z}$ , and  $\hat{w}$ . Moreover, due to the semantical security of the symmetric key and asymmetric key encryption schemes, given ciphertext  $(\hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w})$  the adversary cannot learn anything about the related plaintext, except with a negligible probability,  $\mu(\lambda)$ . Thus, in experiment  $\text{Exp}_3^{A_1}$ , adversary  $\mathcal{A}_1$  cannot tell the value of  $\gamma \in \{0, 1\}$  significantly better than just guessing it, i.e., its success probability is at most  $\frac{1}{2} + \mu(\lambda)$ . Now we move on to property 2, i.e., the privacy of each verdict from  $\mathcal{DR}$ . Due to the privacy-preserving property of the SAP, given  $g_1 \in g$ , a corrupt  $\mathcal{DR}$  cannot learn  $\bar{k}_1$ . So, it cannot find the encryption-decryption key used to generate ciphertext  $\hat{m}, \hat{l}$ , and  $\hat{z}$ . Also, public parameters  $(pk, pp)$  and token  $T_2$  are independent of  $\mathcal{C}$ 's and  $\mathcal{B}$ 's exchanged messages (e.g., payment requests or warning messages) and  $\mathcal{D}_j$ 's verdicts. Furthermore, due to the semantical security of the symmetric key and asymmetric key encryption schemes, given ciphertext  $(\hat{m}, \hat{l}, \hat{z}, \hat{\pi})$  the adversary cannot learn anything about the related plaintext, except with a negligible probability,  $\mu(\lambda)$ . Also, due to the security of the PVE and FVD protocols (i.e., Theorem 2), the adversary cannot link a verdict to a specific auditor with a probability significantly better than the maximum probability,  $Pr'$ , that an auditor sets its verdict to a certain value, i.e., its success probability is at most  $Pr' + \mu(\lambda)$ , even if it is given the final verdicts, except when all auditors' verdicts are 0. We conclude that, excluding the case where all verdicts are 0, given  $(T_2, pk, pp, g, \hat{m}, \hat{l}, \hat{z}, \hat{\pi}, \hat{w}, v)$ , adversary  $\mathcal{A}_3$ 's success probability in experiment  $\text{Exp}_4^{A_2}$  to link a verdict to an auditor is at most  $Pr' + \mu(\lambda)$ .  $\square$

**Theorem 4.** *The PwDR protocol is secure according to Definition 12.*

*Proof.* Due to Lemma 1, the PwDR protocol is secure against a malicious victim. Also, due to lemmas 2 and 3 it is secure against a malicious bank and is privacy-preserving, respectively. Thus, it satisfies all the properties of Definition 12, meaning that the PwDR protocol is indeed secure according to this definition.  $\square$