

Felix Engelmann, Thomas Kerber, Markulf Kohlweiss, and Mikhail Volkhov\*

# Zswap: zk-SNARK Based Non-Interactive Multi-Asset Swaps

**Abstract:** Privacy-oriented cryptocurrencies, like Zcash or Monero, provide fair transaction anonymity and confidentiality, but lack important features compared to fully public systems, like Ethereum. Specifically, supporting assets of multiple types and providing a mechanism to atomically exchange them, which is critical for e.g. decentralized finance (DeFi), is challenging in the private setting. By combining insights and security properties from Zcash and SwapCT (PETS 21, an atomic swap system for Monero), we present a simple zk-SNARKs based transaction scheme, called Zswap, which is carefully malleable to allow the merging of transactions, while preserving anonymity. Our protocol enables multiple assets and atomic exchanges by making use of sparse homomorphic commitments with aggregated open randomness, together with Zcash friendly simulation-extractable non-interactive zero-knowledge (NIZK) proofs. This results in a provably secure privacy-preserving transaction protocol, with efficient swaps, and overall performance close to that of existing deployed private cryptocurrencies. It is similar to Zcash Sapling and benefits from existing code-bases and implementation expertise.

**Keywords:** NIZK, Cryptocurrency, Privacy, Multi-Asset, Exchange

DOI 10.56553/popets-2022-0120

Received 2022-02-28; revised 2022-06-15; accepted 2022-06-16.

## 1 Introduction

Cryptocurrencies are experiencing steady growth not only in terms of general popularity, security, and real-world applicability, but also in terms of diversity of fi-

ancial instruments that can be realized with them. Decentralized finance (DeFi [23]), an umbrella term that covers such financial instruments in the cryptocurrency community, is one of the *raison d'être* of Ethereum, a cryptocurrency popularized by the wide applicability of its smart contract toolchain. One of the foundational pieces of DeFi is the ability (of the distributed ledger) to create user-defined tokens, and trade or exchange them directly on-chain. Ethereum, by providing Turing complete smart contracts and the ERC tokens standards (e.g. ERC20 or non-fungible ERC721), allows building tools such as automated exchanges [24], investment platforms<sup>1,2</sup>, bidding platforms, insurance tools, NFT marketplaces, etc. An important limitation of current DeFi solutions is the public nature of these atomic on-chain exchanges.

While privacy-preserving cryptocurrencies, like Monero and Zcash, are undoubtedly practical, and ample academic research on the problem of private currencies is available [1, 4, 6, 12, 13], much less is known about the *private exchange* of assets [8, 11, 14] — hindering DeFi applications. This is unfortunate, as privacy on the blockchain, such as transaction anonymity or transfer amount secrecy, is not only interesting per se (for the end user), but also changes the financial landscape of the ecosystem. This can be advantageous, for example, if restricting the adversarial view prevents certain harmful behaviour that results from gaming the market (e.g. frontrunning, Miner Extractable Value). Consider a miner which observes a buy-order for an asset that is either from a notorious investor or has an exceptionally large transfer amount. The miner can buy that asset itself early and cheaply before the order drives up the price of the asset. Private swaps of assets can mitigate such attacks.

However, while practically attractive, combining privacy with DeFi tools is challenging. Solutions that try to tackle this question, such as private smart contracts [20] and privacy-friendly decentralized exchanges [3], often find it hard to balance practical applicability with the privacy guarantees they provide.

---

**Felix Engelmann:** IT University of Copenhagen, felixresearch@nlogn.org

**Thomas Kerber:** University of Edinburgh, IOHK, thomas.kerber@iohk.io

**Markulf Kohlweiss:** University of Edinburgh, IOHK, markulf.kohlweiss@ed.ac.uk

**\*Corresponding Author: Mikhail Volkhov:** University of Edinburgh, mikhail.volkhov@ed.ac.uk

---

**1** <https://compound.finance/markets>

**2** <https://polygon.market.xyz/>

Flexibility, so desirable for DeFi, is at odds with privacy, since having both generally requires heavier cryptographic primitives, like general NIZKs, or multi-party computation (MPC).

In this work we do not support a full private DeFi solution, but instead focus on the foundational problem of constructing a private cryptocurrency mechanism that has both support for multiple assets, and an embedded functionality to perform non-interactive atomic assets swaps (as well as regular transfers).

An atomic swap is the exchange of different assets between multiple parties. It has to happen atomically such that all participants get their desired output or the transaction is aborted. A classical example is a foreign currency exchange, where a bank sells a foreign currency to a customer who pays in the local currency. There, atomicity is guaranteed by simultaneously handing over the assets.

The atomic swap protocol we suggest enables untrusted parties to merge transactions off-chain. The transactions themselves only reveal a map of the imbalance for each asset where the sum of its inputs is unequal to the sum of its outputs. Hence, a balanced transaction does not reveal any amounts or types. Especially, the mergers can neither deanonymize senders or receivers nor correlate a subsequent spending of an output by its precise amount and type. With such little trust required in mergers, this very basic functionality already allows creating local exchange markets, where users can send exchange offers (as transactions with a negative imbalance for the asked token and a positive imbalance for the offered), and community-selected participants can match them and merge them to then submit to the blockchain. The anonymity of the system is controlled by users: the system allows both private swaps between several parties, who agree on their exchange off-chain, and bigger exchange pools, as just mentioned. In both cases, the only information that the mergers and other users see is the one necessary to match the offers (total value for each unbalanced type), and it is erased as soon as the transaction is balanced and sent to the ledger. For partial merges, any type with an imbalance of zero is dropped. An open group of participants in a pool may provide sufficient liquidity and maintain a public order book, similar to classical exchanges.

An interesting open question is how to integrate Zswap with private smart contracts to support more elaborate private DeFi solutions. A first step in that direction would be to extend a public smart contract system with minting policies for private assets, to support, e.g. the private trading of NFTs. Finally, our mech-

anism is compatible with existing consensus protocols, such as proof-of-work or stake, and can be viewed as a full cryptocurrency.

## 1.1 Technical Overview

With respect to swaps, a major challenge of many existing solutions is that transaction data is not explicitly separated from the transaction signature, which binds inputs from multiple users together. This, often results in the need for (slow) MPC protocols to construct such a signature jointly. An important insight of our work is that the Zcash ecosystem already implement signature separation. It was first introduced by Zcash Sapling [16] to reduce the size of SNARK circuits for large transactions, and was inherited by the corresponding Shielded Assets<sup>3</sup>, or similar MASP<sup>4</sup> multi-asset protocols. Another inspiration of our work is SwapCT, that realizes atomic swaps on Monero [11]. We continue this direction by designing and proving secure an extension of a Zcash Sapling like system that satisfies the required multi-asset and atomic swap properties.

Zcash Sapling, as opposed to the Zerocash paper, separates the validation of inputs and outputs (each input and output requires a separate NIZK) from the transaction balancing, which is done using homomorphic commitments and a Schnorr-like binding signature. The binding signature “seals” the inputs and outputs in place, forbidding adding or removing any extra inputs or outputs. Instead, in our work we use a sparse multi-value Pedersen commitments and relax the signature, allowing transactions to be non-interactively merged together. While inspired by SwapCT, we deal with different challenges specific to the zk-SNARK setting and take advantage of the existing signature separation in Zcash Sapling. The explicit signature separation opens more space for potential transaction malleability, and should be taken with care. To our knowledge our work is the first to extend the security analysis of [4] in that direction.

Our modelling relies on a one-time account (OTA) scheme to anonymously and confidentially create and store value in notes. This is an abstraction from existing protocols in Zcash and Monero. The spending of input

<sup>3</sup> Previously “user-defined assets” (UDA) or UIT, see ZIP 220: <https://github.com/zcash/zcash/issues/830>

<sup>4</sup> <https://github.com/anoma/masp/blob/main/docs/multi-asset-shielded-pool.pdf>

notes, stored in a Merkle tree, is authorized by a simulation extractable (SE) non-interactive zero-knowledge proof (NIZK). Double spends are prevented by proving the correctness of a deterministic nullifier, marking an input as spent. To connect inputs to newly created output notes, we use sparse homomorphic commitments to value-type pairs, which can be summed to check that the transaction is well-balanced. Importantly, the values and types remain hidden even when we publish their aggregated randomness, which is necessary to facilitate transaction merging.

Our protocol bears similarities to variants of Zcash and SwapCT which we compare in more detail:

**Zcash Sapling:** Compared to Sapling, we add multi-asset support and remove the authorization signature which is replaced by the SE NIZK. We also do not need a binding signature over all intermediate commitments — instead we directly publish the randomness. This opens room for controlled malleability, and enables non-interactive joint transaction generation (by merging) without MPC, which can be used to implement swaps mechanics. The multi-asset aspect of our scheme is very similar to Shielded Assets or MASP, both of which extend Sapling but do not provide a mechanism for atomic swaps. We also provide a rigorous theoretical formalization, which is lacking for both Sapling and its Shielded Asset and MASP variants, and can be of independent interest.

Conversely, the technical discussions and implementation effort that went into the Shielded Assets ZIP 220 and the MASP implementation are valuable starting points for the deployment of Zswap as part of a larger Zcash like blockchain. In particular, deployment of our new transactions into Zcash requires a hard fork and creates a new shielded pool holding typed notes. Transferring tokens from an existing shielded pool is possible by explicitly adding the Zcash type to the notes.

**SwapCT:** Our work is inspired by SwapCT which also provides atomic swaps and transaction merging, but on top of Monero (ring based anonymity). We simplified their scheme to unify their offers and transactions. With the use of SNARKs (requiring a trusted or transparent setup) our simpler construction no longer needs their anonymously aggregatable signatures. We achieve better confidentiality for unfinished transactions as we only reveal a total imbalance instead of per input and output values. Another important difference is that the use of SNARKS allows to efficiently support large rings. This motivates the use of a single global ring in our formalization.

**Formalization:** We formalize our protocol using game based definitions, which generalize Zerocash definitions, but are also compatible with Monero except for using a single global ring (as in Zcash and differing from SwapCT).

Zerocash [4] formalizes *Ledger indistinguishability*, *Transaction non-malleability* and *Balance*. Our balance definition is very similar with the addition of checking multiple assets. We replace indistinguishability with the privacy definition, as our transactions may be circulated off-chain. The non-malleability of merged transactions is inspired by Zerocash and SwapCT [11]. Instead of the strong non-malleability, it provides a relaxed property that controls malleability (transaction merging) and only requires theft prevention, namely that an adversary cannot steal from the intended recipients of a transaction, or split unbalanced transactions, rerouting honest funds. Modelling this malleability relies on a variant of hiding of the Pedersen commitment scheme, when the joint commitment randomness is revealed, which we call HID-OR<sup>5</sup>. This property does not require any additional security assumptions.

## 1.2 Our Contributions

In this work we present the following results:

- We specify a formal model for a multi-asset Zcash system with swaps that builds on top of a One-Time Account (OTA) System, abstracting a nullifier-like private UTXO mechanism. The OTA model and our proof techniques could be of independent interest for proving systems such as Zcash and Monero secure.
- We provide a minimal practical instantiation of non-interactive private atomic-swaps. It is based on a simplified version of Zcash that removes authorization and blinding signatures.
- We prove our construction secure under commonly used assumptions similar to the ones used in Zerocash. This validates the removal of the Zcash signatures and shows that the perfect hiding and binding properties of spend and output commitments is sufficient for security.
- We implement and evaluate our protocol: our merging mechanism is very effective, and all performance overheads of our construction, as compared to the basic single-asset protocol without swaps, are small.

---

<sup>5</sup> HIDing with Open Randomness

### 1.3 Related Work

Surprisingly little academic work *directly* addresses exchange of multiple assets in a private ledger. Undoubtedly, such a functionality can be achieved through certain generic private smart contract solutions [17, 18, 22], but their flexibility comes at a cost of non-negligible performance overhead, since they often require heavy primitives like SNARKs over big contract code-dependant (or even universal) circuits. The performance of universal zero-knowledge based constructions such as ZEXE [5] requires minutes of proving time for ten times the constraints compared to our one second prover runtime. This is why we would like to consider systems with such a functionality embedded directly.

Several solutions take the route of extending the vanilla Zerocash. Ding et al. [10] propose a solution supporting multiple assets, but with no exchange mechanism, and with public asset types. Gao et al. [14] construct a transaction system specifically for exchanging assets which is based on storing debt in sibling notes which are only spendable if the debt is settled. Its inefficiency results from requiring multiple persisted transactions per swap.

On the other hand, Confidential Assets [21] uses a commitment construction that hashes an asset type descriptor to the bases of an extended Pedersen commitment such that the resulting commitments are additively homomorphic, which facilitates proving the balancing of amounts. We use this sparse commitment scheme as part of our construction. A similar solution by Zheng et al. [25] exists for the Mumblewimble private cryptocurrency. Both these works do not provide sender and receiver anonymity.

Sparse homomorphic commitments have the drawback that they can only store type-amount pairs. A more flexible approach is to use a hash-based commitment scheme for notes to store a vector of attributes. This is used in Zcash’s multi-asset ZIP 220, Shielded Assets, or similarly MASP. To achieve balancing, these protocols prove equality between the type-value pair of a note and a sparse homomorphic commitment. We follow the same approach. Another system that *does not* rely on homomorphic commitments is Stellar<sup>6</sup>, which instead uses shuffle proofs.

Finally, some works emphasize the exchange and offer matching functionality. Manta [8] describes a

privacy-preserving decentralized exchange (DEX) based on an automated market maker (AMM) scheme which works without a second party but does not hide the types, which is an inherent limitation of the AMM approach. Another idea is to privately exchange assets between different systems in cross-chain atomic swaps [9]. In contrast to both, we propose a more basic mechanism within a single blockchain, and leave it open to implement the concrete offer matching algorithm on top of Zswap. This allows, and will likely enable more powerful DeFi applications, since Zswap provides a more flexible interface to the application layer.

## 2 Preliminaries

Let us introduce our notation. We write PPT for probabilistic polynomial-time. All our sets are multisets by default. The set minus operation  $A \setminus B$  removes as many values in  $A$  as there are in  $B$ : if  $B$  has more values of the same type than in  $A$ , the number of elements in the resulting set is 0. The function  $\text{ToSet}(S)$  removes all duplicate entries in the multiset  $S$  and sorts the result according to some predefined ordering, essentially converting  $S$  to a “proper” set. We use  $[n] := 1 \dots n$ , and  $|S|$  for the cardinality of the set  $S$ ,  $|\vec{x}|$  is the length if  $\vec{x}$  is a vector, thus  $[|S|]$  means  $1 \dots |S|$ .

As of pseudocode conventions, abortion by default means returning  $\perp$  immediately, and assertions abort on failure. When the result of an assignment cannot be properly pattern-matched, e.g.  $(a, b) \leftarrow \text{foo}$  (function returns  $\perp$  or a different type), the execution aborts. The symbol  $@$  denotes temporary variable assignment in the pattern-matching cases, e.g.  $(a@(x, y), z) \leftarrow \text{Foo}()$  means  $((x, y), z) \leftarrow \text{Foo}(); a \leftarrow (x, y)$ .

### 2.1 Sparse Homomorphic Commitments

Our first building block is a commitment scheme which is additively homomorphic for an exponential number of possible domains, also called types, however each commitment is sparsely populated, i.e. few domains have a non-zero value. Between domains, there exists no homomorphic property. Constructions were proposed by [21] and formalized by [7] which resemble vector Pedersen commitments with ad-hoc generators.

<sup>6</sup> <https://github.com/stellar/slingshot/blob/main/spacesuit/spec.md>

**Definition 2.1** (Sparse Homomorphic Commitment).

A SHC scheme consists of the algorithms  $\text{ComSetup}$  and  $\text{Commit}$  defined as follows:

$\mathfrak{p} \leftarrow \text{ComSetup}(1^\lambda)$  takes the security parameter  $\lambda$  and outputs the public parameters  $\mathfrak{p}$  implicitly provided to  $\text{Commit}$ . It specifies a type space  $\mathbb{T}$   
 $\text{com} \leftarrow \text{Commit}(\{\{\text{ty}_i, a_i\}_{i=1}^n, \text{rc}\})$  takes a set of distinct types  $\text{ty}_i$  and values  $a_i$ , and a randomness  $\text{rc}$ , and outputs a commitment  $\text{com}$ . We write  $\text{Commit}(\text{ty}, a, \text{rc})$  for  $\text{Commit}(\{\{\text{ty}, a\}, \text{rc}\})$ .

**Definition 2.2** (Homomorphism). A SHC scheme is homomorphic, if there exists an efficient operation  $\oplus$ , such that for any  $\text{ty} \in \mathbb{T}$  it holds that

$$\text{Commit}(\{\{\text{ty}_i, a_i\}_1^n, \text{rc}\}) \oplus \text{Commit}(\{\{\text{ty}_i, a'_i\}_1^n, \text{rc}'\}) = \text{Commit}(\{\{\text{ty}_i, a_i + a'_i\}_1^n, \phi(\text{rc}, \text{rc}')\})$$

for some function  $\phi$ , while zero values do not affect the input set:  $\text{Commit}(\{\{\text{ty}, 0\}, \text{rc}\}) = \text{Commit}(\emptyset, \text{rc})$ .

In other words, the commitment is additive on the same type, but acts like a vector commitment on distinct types.

We require standard binding and hiding properties.

**Definition 2.3** (Commitment Binding). For any  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \leftarrow \text{Setup}(1^\lambda)$  and any PPT adversary  $\mathcal{A}$ , it holds that

$$\Pr \left[ \begin{array}{l} \{\{\text{ty}_b, a_b, \text{rc}_b\}_{b=0}^1\} \leftarrow \mathcal{A}(\mathfrak{p}); \\ \text{com}_0 = \text{Commit}(\{\{\text{ty}_0, a_0\}, \text{rc}_0\}); \\ \text{com}_1 = \text{Commit}(\{\{\text{ty}_1, a_1\}, \text{rc}_1\}); \\ \text{com}_0 = \text{com}_1 \wedge \\ (\text{ty}_0, a_0, \text{rc}_0) \neq (\text{ty}_1, a_1, \text{rc}_1) \end{array} \right] \leq \text{negl}(\lambda)$$

**Definition 2.4** (Commitment Hiding). For any  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \leftarrow \text{Setup}(1^\lambda)$  and any stateful PPT adversary  $\mathcal{A}$ , it holds that

$$\Pr \left[ \begin{array}{l} \{\{\text{ty}_i, a_i\}_{i=0}^1\} \leftarrow \mathcal{A}(\mathfrak{p}); \\ b \xleftarrow{\$} \{0, 1\}, \text{rc} \xleftarrow{\$} \mathbb{R}; \\ \text{com} \leftarrow \text{Commit}(\{\{\text{ty}_b, a_b\}, \text{rc}\}); \\ b' \leftarrow \mathcal{A}(\text{com}) : \\ b' = b \end{array} \right] \leq \frac{1}{2} + \text{negl}(\lambda)$$

To instantiate a SHC we use a well-known Pedersen-based construction derived from the [21], which we will simply call “sparse Pedersen commitment scheme”. Let  $\mathbb{G}$  be a cyclic group of order  $q$  with a generator  $G$  in which the discrete logarithm assumption holds. Additionally, we use a cryptographic hash function  $H : \mathbb{T} \rightarrow \mathbb{G}$ . Then  $\text{Commit}(\{\{\text{ty}_i, a_i\}, \text{rc}\}) := (\prod H(\text{ty}_i)^{a_i}) G^{\text{rc}}$  with  $a_i, \text{rc} \in \mathbb{Z}_q$ .

$\text{HID-OR}_{\mathcal{A}}^b(1^\lambda)$

---

$\mathfrak{p} \leftarrow \text{Setup}(1^\lambda)$   
 $\{\{\text{ty}_t, a_t\}, \{\text{ty}'_t, a'_t\}\}_{t=0}^1 \leftarrow \mathcal{A}(\mathfrak{p})$   
 $\Delta_{\text{ty}, t} := \{a_t \text{ if } \text{ty}_t = \text{ty} \text{ else } 0\} - \{a'_t \text{ if } \text{ty}'_t = \text{ty} \text{ else } 0\}$   
 assert  $\forall \text{ty} \in \bigcup \text{ty}_t \cup \bigcup \text{ty}'_t : \Delta_{\text{ty}, 0} = \Delta_{\text{ty}, 1}$   
 $r, r' \xleftarrow{\$} \mathbb{R}$   
 $\text{com} \leftarrow \text{Commit}(\{\{\text{ty}_b, a_b\}, r\}); \text{com}' \leftarrow \text{Commit}(\{\{\text{ty}'_b, a'_b\}, r'\})$   
 $b' \leftarrow \mathcal{A}(\text{com}, \text{com}', \text{rc} = r - r')$   
 return  $b'$

Fig. 1. Hiding with Open Randomness Game

Finally, we observe the following property of the sparse Pedersen commitment scheme that guarantees that the adversary cannot distinguish between two pairs of commitments that sum to the same values (per type), even if we reveal their common randomness. This is the main property used in our swaps: we will use it to argue that transaction can be merged (to join two unbalanced swap offers), but not split apart (so swap offers cannot be adversarially modified).

**Lemma 2.1** (Hiding with Open Randomness). The sparse Pedersen commitment scheme is perfectly hiding with open randomness in the RO model. By this we mean that if for all  $\lambda \in \mathbb{N}$  and all PPT  $\mathcal{A}$ :

$$\Pr[\text{HID-OR}_{\mathcal{A}}^1(1^\lambda) = 1] - \Pr[\text{HID-OR}_{\mathcal{A}}^0(1^\lambda) = 0] = 0$$

where the game is defined in Fig. 1.

The proof of this statement is presented in Appendix C and based on perfect hiding together with the equivocation of the Pedersen commitment scheme.

## 2.2 Proofs and Signatures of Knowledge

A non-interactive zero-knowledge proof of knowledge (or just NIZK) for the language  $\mathcal{L}$  is a proving system consisting of algorithms (Setup, Prove, Verify, Sim). The first algorithm allows to create a common reference string (CRS),  $\text{crs} \leftarrow \text{Setup}(1^\lambda)$  that is used as an input to Prove and Verify. To obtain a non-interactive proof that  $\text{stmt} \in \mathcal{L}$  with the corresponding witness  $w$  (s.t.  $(\text{stmt}, w) \in \mathcal{R}_{\mathcal{L}}$ ) one executes  $\pi \leftarrow \text{Prove}(\text{crs}, \text{stmt}, w)$ . To then verify this prove one runs, correspondingly  $\text{Verify}(\text{crs}, \text{stmt}, \pi)$ , which returns 0 or 1. We provide the standard definitions and security properties in Appendix A.

Signatures of Knowledge are NIZK proofs which are bound to a specific message  $m$ . For a simulation extractable (SE) NIZK protocol, this is done by including

$m$  into the statement without asserting anything about  $m$  in the relation. Still, the SE property invalidates the proof, if the statement is changed.

### 3 One-Time-Account Scheme

To highlight our novel transaction mechanism, we first describe the *one-time account (OTA) scheme* that we use to model the mechanics of underlying accounts our transactions use. The OTA scheme may be seen as an anonymous version of an unspent transaction output (UTXO) system (e.g. the one used in Bitcoin); it generalises accounts of privacy-preserving transaction systems such as Zcash or Monero. Instead of creating transaction outputs including a long term identity as with plain UTXO, an OTA scheme generates a unique, anonymous, one-time account for each transaction output. To support the UTXO functionality to decide if a one-time-account is still valid, the OTA scheme allows generating a unique *nullifier* which anonymously marks it as spent. A duplicate nullifier indicates that the same one-time account is used.

The OTA scheme is used as follows. After a system *Setup*, each participant generates their credentials with *KeyGen*. Anyone with the public key can then non-interactively derive a one-time account, *also called “note”*, with *Gen*. Each note contains a set vector of attributes. Most importantly, one attribute is the amount the note represents. In our multi-type setting, a second attribute is the type of the amount. Other systems may make use of additional note attributes such as e.g. time locks — a timestamp identifier, allowing a note to be spendable only after a specified block number. To allow the intended recipient to recover the note, it is accompanied by an encryption  $C$  calculated by *Enc*. With the secret key, note and ciphertext, the original owner *Receives* the note and is then able to create the corresponding nullifier with *NulEval*, a unique serial value characterizing the note that cannot be predicted by anyone else than the owner.

We emphasize that any party can create an OTA account for anyone knowing only their public key. The OTA itself does not allow any claims to the value stored inside. In our system, only the OTAs which are included as outputs of a valid, balanced and persisted transaction can be claimed by their owners in subsequent transactions. The transaction then enforces that sufficient inputs were consumed to cover for the output. To mint coins an OTA needs to be included on the ledger in an

unbalanced transaction for the newly minted type. This transaction needs to be accepted by the ledger rules. This behaviour might as well be governed by a smart contract.

**Definition 3.1.** *A One-Time-Account (OTA) scheme consists of the five PPT algorithms (Setup, KeyGen, Gen, Receive, NulEval) defined as follows:*

$\mathbf{p} \leftarrow \text{Setup}(1^\lambda)$ : *takes the security parameter  $\lambda$  and outputs the public parameters  $\mathbf{p}$  which are implicitly provided to the subsequent algorithms. This includes the note randomness space  $\mathbb{S}$ , a message space  $\mathbb{M}$ , and an encryption randomness space  $\Xi$ .*

$(\mathbf{sk}, \mathbf{pk}) \leftarrow \text{KeyGen}(1^\lambda)$ : *generates a key pair  $(\mathbf{sk}, \mathbf{pk})$ . We assume a function  $P(\mathbf{sk}) = \mathbf{pk}$  for generating the public key from a secret key.*

$\text{note} \leftarrow \text{Gen}(\mathbf{pk}, \vec{a}, r)$ : *takes a public key  $\mathbf{pk}$ , a vector of attributes  $\vec{a} \in \mathbb{M}^{|\vec{a}|}$ , where, by agreement, the first might be an amount, and randomness  $r \in \mathbb{S}$ . It outputs a one-time-account, known as note.*

$C \leftarrow \text{Enc}(\mathbf{pk}, (\vec{a}, r), \xi)$ : *encrypts the attributes  $\vec{a}$  and the randomness  $r$  to the public key  $\mathbf{pk}$  with additional randomness  $\xi$ . It outputs a ciphertext  $C$ .*

$(\vec{a}, r)/\perp \leftarrow \text{Receive}(\text{note}, C, \mathbf{sk})$ : *if the note and ciphertext  $C$  belongs to the secret key  $\mathbf{sk}$ , the algorithm outputs the vector of attributes  $\vec{a}$  and the randomness  $r$  or fails otherwise.*

$\text{nul} \leftarrow \text{NulEval}(\mathbf{sk}, r)$ : *Takes a secret key  $\mathbf{sk}$  and a randomness  $r$  and outputs a nullifier  $\text{nul}$ .*

*In addition, the algorithms Gen and NulEval must be efficiently provable in zero-knowledge. More formally, the construction must provide a NIZK protocol for the following languages:*

$$\begin{aligned} \mathcal{L}^{\text{nul}} &= \{(\text{note}, \text{nul}) \mid \exists(\mathbf{sk}, \vec{a}, r) : \\ &\quad \text{note} = \text{Gen}(P(\mathbf{sk}), \vec{a}, r) \wedge \text{nul} = \text{NulEval}(\mathbf{sk}, r)\} \\ \mathcal{L}^{\text{open}} &= \{\text{note} \mid \exists(\mathbf{pk}, \vec{a}, r) : \text{note} = \text{Gen}(\mathbf{pk}, \vec{a}, r)\} \end{aligned}$$

*The language  $\mathcal{L}^{\text{open}}$  may optionally be extended, such that elements of  $\vec{a}$  have relations to other commitments.*

**Definition 3.2** (OTA Correctness). *An OTA scheme is correct if for any  $\lambda \in \mathbb{N}$  with  $\mathbf{p} \in \text{Setup}(1^\lambda)$  it holds that any honestly generated note is receivable. Formally, for every  $(\mathbf{sk}, \mathbf{pk}) \in \text{KeyGen}(\mathbf{p})$ , every  $(\vec{a}, r) \in \mathbb{M}^{|\vec{a}|} \times \mathbb{S}$  and  $\xi \in \Xi$  it holds that  $\text{Receive}(\text{Gen}(\mathbf{pk}, \vec{a}, r), \text{Enc}(\mathbf{pk}, (\vec{a}, r), \xi), \mathbf{sk}) = (\vec{a}, r)$*

Once a note is created, it must bind the attributes and prevent opening the note to a different vector, even for the owner with the correct secret key.

**Definition 3.3** (OTA Binding). *An OTA scheme is binding with regard to the accounts created and the vector of attributes if for any  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \in \text{Setup}(1^\lambda)$  and any PPT adversary  $\mathcal{A}$ , it holds that*

$$\Pr \left[ \begin{array}{l} (\mathfrak{pk}_0, r_0, \vec{a}_0, \mathfrak{pk}_1, r_1, \vec{a}_1) \leftarrow \mathcal{A}(\mathfrak{p}) \\ \text{note}_0 \leftarrow \text{Gen}(\mathfrak{pk}_0, \vec{a}_0, r_0) \\ \text{note}_1 \leftarrow \text{Gen}(\mathfrak{pk}_1, \vec{a}_1, r_1) : \\ \text{note}_0 = \text{note}_1 \wedge \vec{a}_0 \neq \vec{a}_1 \end{array} \right] \leq \text{negl}(\lambda)$$

The following privacy property assures that a note and its ciphertext do not leak who the note belongs to and what attributes it stores.

**Definition 3.4** (OTA Privacy). *Consider the following oracle:*

$$\begin{array}{l} \mathcal{O}_{\text{Rcv}}^{\text{note}^*, C^*}(i, \text{note}, C) : \\ \quad \text{assert } \text{note} \neq \text{note}^* \wedge C^* \neq C \\ \quad \text{return } \text{Receive}(\text{note}, C, \text{sk}_i) \end{array}$$

*An OTA scheme is private if for any  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \in \text{Setup}(1^\lambda)$  and any stateful PPT adversary  $\mathcal{A}$ , it holds that*

$$\Pr \left[ \begin{array}{l} (\mathfrak{sk}_0, \mathfrak{pk}_0) \xleftarrow{\$} \text{KeyGen}() \\ (\mathfrak{sk}_1, \mathfrak{pk}_1) \xleftarrow{\$} \text{KeyGen}() \\ (i_0, \vec{a}_0, i_1, \vec{a}_1) \leftarrow \\ \quad \mathcal{A}^{\mathcal{O}_{\text{Rcv}}^{\perp, \perp}}(\mathfrak{p}, \mathfrak{pk}_0, \mathfrak{pk}_1) \\ b \xleftarrow{\$} \{0, 1\}, r \xleftarrow{\$} \mathbb{S}, \xi \xleftarrow{\$} \Xi \\ \text{note}^* \leftarrow \text{Gen}(\mathfrak{pk}_{i_b}, \vec{a}_b, r) \\ C^* \leftarrow \text{Enc}(\mathfrak{pk}_{i_b}, (\vec{a}_b, r), \xi) \\ b' \leftarrow \mathcal{A}^{\mathcal{O}_{\text{Rcv}}^{\text{note}^*, C^*}}(\text{note}^*, C^*) : \\ b = b' \wedge |\vec{a}_0| = |\vec{a}_1| \end{array} \right] \leq \frac{1}{2} + \text{negl}(\lambda)$$

The privacy game implicitly subsumes (1) note and ciphertext hiding, and (2) note and encryption anonymity (key privacy). In the first case,  $\mathcal{A}$  cannot decide the content of the note or the ciphertext; in the second, it cannot decide which key was used to create it. E.g. if note is not hiding,  $\mathcal{A}$  efficiently wins the game by first returning two different vectors  $\vec{a}_0 \neq \vec{a}_1$  and then distinguishing the note  $\text{note}_b$  according to the attribute vector.

We require notes to be *unique* when generated with honest randomness:

**Definition 3.5** (Note Uniqueness). *A binding and private OTA scheme satisfies honestly generated note uniqueness: for all PPT  $\mathcal{A}$ ,*

$$\Pr \left[ \begin{array}{l} ((\mathfrak{pk}_0, \vec{a}_0), (\mathfrak{pk}_1, \vec{a}_1)) \leftarrow \mathcal{A}(1^\lambda) \\ r_0, r_1 \xleftarrow{\$} \mathbb{S} : \\ \text{Gen}(\mathfrak{pk}_0, \vec{a}_0, r_0) = \text{Gen}(\mathfrak{pk}_1, \vec{a}_1, r_1) \end{array} \right] \leq \text{negl}(\lambda)$$

The next property, similar to the tagging scheme of Omniring [19], captures the requirement that nullifiers produced with  $\text{sk}$  must be only “predictable” for the party that holds  $\text{sk}$ . This is achieved by requiring  $\text{NulEval}$  to behave like a pseudorandom function.

**Definition 3.6** (Nullifier Pseudorandomness). *Let  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \in \text{Setup}(1^\lambda)$ ,  $(\text{sk}, \text{pk}) \xleftarrow{\$} \text{KeyGen}()$ , and  $f$  be randomly sampled function on the range  $\{0, 1\}^{|\mathbb{S}|} \rightarrow \{0, 1\}^{|\text{NulEval}(\text{sk}, \cdot)|}$ . An OTA scheme nullifier is pseudorandom if for any PPT adversary  $\mathcal{A}$ , it holds that*

$$\frac{\Pr[\mathcal{A}^{\text{Receive}(\cdot, \cdot, \text{sk}), f(\cdot)}(\text{pk}) = 1] - \Pr[\mathcal{A}^{\text{Receive}(\cdot, \cdot, \text{sk}), \text{NulEval}(\text{sk}, \cdot)}(\text{pk}) = 1]}{\Pr[\mathcal{A}^{\text{Receive}(\cdot, \cdot, \text{sk}), \text{NulEval}(\text{sk}, \cdot)}(\text{pk}) = 1]} \leq \text{negl}(\lambda)$$

Finally, to detect duplicate use of the same note, each is assigned a unique nullifier. Even with knowledge of the secret key, it is not possible to create two different nullifiers for the same note. The separate secret keys are important for constructions based on algebraic nullifiers. E.g. Omniring creates tags as  $g^{\frac{1}{\text{sk}+r}}$  for a generator  $g$ , so randomness may be “traded” for secret key. Allowing only one note with a single secret key would not capture the realistic setting where an adversary controls multiple correlated accounts.

**Definition 3.7** (Nullifier Uniqueness). *An OTA scheme satisfies nullifier uniqueness if for any  $\lambda \in \mathbb{N}$  with  $\mathfrak{p} \in \text{Setup}(1^\lambda)$  and any PPT adversary  $\mathcal{A}$ , it holds that*

$$\Pr \left[ \begin{array}{l} (\mathfrak{sk}_0, r_0, \vec{a}_0, \mathfrak{sk}_1, r_1, \vec{a}_1) \leftarrow \mathcal{A}(\mathfrak{p}); \\ \text{note}_0 \leftarrow \text{Gen}(P(\mathfrak{sk}_0), \vec{a}_0, r_0); \\ \text{note}_1 \leftarrow \text{Gen}(P(\mathfrak{sk}_1), \vec{a}_1, r_1) : \\ \text{note}_0 = \text{note}_1 \wedge \\ \text{NulEval}(\mathfrak{sk}_0, r_0) \neq \text{NulEval}(\mathfrak{sk}_1, r_1) \end{array} \right] \leq \text{negl}(\lambda)$$

## 4 Zswap Scheme

A Zswap scheme is an extension of an OTA scheme which allows creating ( $\text{SignTx}$ ), merging ( $\text{MergeSig}$ ), and verifying ( $\text{Verify}$ ) transactions that transfer coins between OTA accounts.  $\text{SignTx}$  takes a pre-transaction  $\text{ptx}$  as input and produces a transaction signature  $\sigma$ ,  $\text{MergeSig}$  combines transaction signatures  $\sigma_1, \dots, \sigma_n$ , and  $\text{Verify}$  verifies a signature  $\sigma$  for a transaction  $\text{tx}$ . A signature  $\sigma$  is viewed separately from its transaction  $\text{tx}$  (created by  $\text{tx} \leftarrow \text{CompleteTx}(\text{ptx})$  defined below), and not contained in it.

We first present auxiliary algorithms for Zswap in Fig. 2 for creating pre-transactions and transactions. These use the OTA scheme algorithms only in a black-box manner and simplify the exposition when defining the security properties. The first set of functions glues the Zswap and OTA schemes together.

- BuildPTx constructs inputs for SignTx from  $I, O$  instructions and a set of secret keys SK. Inputs  $I$  consist of a list of existing (note,  $C$ ) notes, outputs  $O$  consist of a list of public key, value, and type triplets  $(pk^T, a^T, ty^T)$  for creating the output notes. By receiving input notes and generating output notes it produces a pre-transaction information ptx.
- CompleteTx constructs a transaction tx for Verify from a pre-transaction ptx similar to SignTx.
- MergeTx creates a new transaction tx by combining a set of existing transactions  $tx_1, \dots, tx_n$ . This is the non-cryptographic analogue of MergeSig.
- CheckPTx(ptx, st) checks whether pre-transaction is valid w.r.t. st, which is used in the correctness property. It is easy to see that for  $ptx = \text{BuildPTx}(\cdot, \cdot, st, \cdot)$  we have  $\text{CheckPTx}(ptx, st) = 1$ .
- TryReceive attempts to “receive” a note note by decrypting its ciphertext  $C$  using any of the set SK of available secret keys: if an input can be received with one of the keys it also computes the note’s nullifier.

**Definition 4.1.** A Zswap transaction scheme, built on top of an OTA scheme, consists of a tuple of PPT algorithms (Setup, SignTx, MergeSig, Verify) defined as follows:

- $\rho \leftarrow \text{Setup}(1^\lambda)$  takes the security parameter  $\lambda$  and outputs public parameters  $\rho$  which are implicitly given to all the following algorithms. Setup is called once when a Zswap system is initialized.
- $\sigma \leftarrow \text{SignTx}(ptx, st)$  takes a pre-transaction  $ptx = (\mathcal{S}, \mathcal{T})$  where
  - st is the current state of valid previously issued notes. st.MT is the Merkle tree containing notes as leaves, and st.NF is the set of used nullifiers. If a state st is the evolution of a previous state  $st'$ , we define a function  $\text{ev}(st', st) := 1$  and 0 otherwise.
  - $\mathcal{S} = \{(sk_i^S, \text{note}_i^S, \text{nul}_i, \text{path}_i, (a_i^S, ty_i^S), r_i^S)\}_{i=1}^{|\mathcal{T}|}$  is a set of inputs with a nullifier  $\text{nul}_i$  corresponding to the note  $\text{note}_i^S$ , and stored in the current state at the given path  $\text{st.MT}[\text{path}_i]$ , secret key  $sk_i^S$ , amount  $a_i^S$  and type  $ty_i^S$  with input OTA notes’ randomness  $r_i^S$ .

- $\mathcal{T} = \{(pk_i^T, \text{note}_i^T, (a_i^T, ty_i^T), r_i^T)\}_{i=1}^{|\mathcal{T}|}$  is a set of (output) notes  $\text{note}_i^T$  with amount  $a_i^T$ , type  $ty_i^T$  and output OTA notes’ randomness  $r_i^T$ .

It outputs a signature  $\sigma$  as authorization to spend the inputs  $\mathcal{S}$  on the given outputs  $\mathcal{T}$ .

$b \leftarrow \text{Verify}(st, tx, \sigma)$  takes a transaction tx, a signature  $\sigma$  and returns a bit  $b$  representing the validity of the transaction w.r.t. the state st.

$\sigma \leftarrow \text{MergeSig}(\{\sigma_i\}_{i=1}^n)$  takes  $n$  transaction signatures and generates a combined signature  $\sigma$  valid for the union of the transactions. To merge the corresponding transactions  $tx_1, \dots, tx_n$  together, we use the function MergeTx defined in Fig. 2. It just concatenates their input nullifiers, output notes, and sums their  $\Delta_{ty}$  for each ty. As of the st, they must be equal in the transactions that need to be merged, but we can assume that st is only updated once an epoch, which is set to be a time interval long enough for transactions to be merged.

## 4.1 Atomic Swap Example

After presenting all necessary algorithms, we show a small example of how they interact to create an atomic swap transaction between two parties, Alice and Bob. First, the system is created by Setup. Each participant joining, generates their key pair (sk, pk) with KeyGen. The creation of new assets is delegated to an external consensus mechanism which updates the global state of the system according to an agreed policy. At one point, Alice and Bob have to be the beneficiary of a transaction. They notice this by calling TryReceive with their secret key sk on every published transaction output (note,  $C$ ). If they successfully receive a (note,  $C$ ) to an amount  $a$  and type ty, they keep it for when they want to spend it. Let’s assume Alice received a note of 10\$ and Bob has a note of 10€. Alice now wants euros and Bob dollars. They assume an exchange rate of 7:5 and proceed to generate their pre-transactions  $ptx_i$ . Each party calls BuildPTx with their note as input instruction  $I$ . As output instructions  $O$ , Alice sends 5€ to her key pk. To make sure that someone fulfils the offer, she creates a change output to herself with only 2\$, leaving 1\$ as incentive. Bob performs the same. He inputs the 10€ note and generates outputs of 7\$ and 5€ to himself.

The resulting pre-transactions are signed by both parties respectively with SignTx to get signatures  $\sigma_i$ . The final transactions  $tx_i$  are generated from the pre-transactions  $ptx_i$  by CompleteTx and can be verified against their signature  $\sigma_i$  with Verify. Note that both



<pre> BuildPTx(<math>I, O, \text{st}, \text{SK}</math>) <math>\mathcal{S}, \mathcal{T} \leftarrow \emptyset</math> for (note<math>^S, C^S</math>) <math>\in I</math> do   path <math>\leftarrow \text{st.MT.getPath}(\text{note}^S)</math>   assert path <math>\neq \perp</math>   (sk<math>^S, \text{nul}, (a^S, \text{ty}^S), r^S) \leftarrow \text{TryReceive}(\text{note}, C, \text{SK})</math>   <math>\mathcal{S} := \mathcal{S} \cup (\text{sk}^S, \text{note}^S, \text{nul}, \text{path}, (a^S, \text{ty}^S), r^S)</math>   assert {nuli} are distinct and <math>\forall i. \text{nul}_i \notin \text{st.Nf}</math> for (pk<math>^T, a^T, \text{ty}^T</math>) <math>\in O</math> do <math>r_i^T \xleftarrow{\\$} \mathbb{S}, \xi \xleftarrow{\\$} \Xi</math>   note<math>^T \leftarrow \text{OTA.Gen}(\text{pk}^T, (a^T, \text{ty}^T), r^T)</math>   <math>C^T \leftarrow \text{OTA.Enc}(\text{pk}^T, ((a^T, \text{ty}^T), r^T), \xi)</math>   <math>\mathcal{T} := \mathcal{T} \cup (\text{pk}^T, \text{note}^T, C^T, (a^T, \text{ty}^T), r^T)</math> return ptx@(<math>\mathcal{S}, \mathcal{T}</math>)  CheckPTx(ptx@(<math>\mathcal{S}, \mathcal{T}</math>), st) Parse <math>\mathcal{S}</math> as <math>\{(\text{sk}_i^S, \text{note}_i^S, \text{nul}_i, \text{path}_i, (a_i^S, \text{ty}_i^S), r_i^S)\}_{i=1}^{ \mathcal{S} }</math> Parse <math>\mathcal{T}</math> as <math>\{(\text{pk}_i^T, \text{note}_i^T, C_i^T, (a_i^T, \text{ty}_i^T), r_i^T)\}_{i=1}^{ \mathcal{T} }</math> % <math>C_i^T</math> is left to the receiver to verify assert {nuli} are distinct and <math>\forall i. \text{nul}_i \notin \text{st.Nf}</math> for <math>i \in [ \mathcal{S} ]</math> do   assert st.MT[path<math>_i</math>] = note<math>_i^S</math>   assert nul<math>_i = \text{OTA.NulEval}(\text{sk}_i^S, r_i^S)</math>   assert note<math>_i^S = \text{OTA.Gen}(\text{OTA.P}(\text{sk}_i^S), (a_i^S, \text{ty}_i^S), r_i^S)</math> assert <math>\forall i \in [ \mathcal{T} ] : \text{note}_i^T = \text{OTA.Gen}(\text{pk}_i^T, (a_i^T, \text{ty}_i^T), r_i^T)</math> return 1 </pre>	<pre> CompleteTx(ptx@(<math>\mathcal{S}, \mathcal{T}</math>)) Parse <math>\mathcal{S}</math> as <math>\{(\cdot, \cdot, \text{nul}_i, \cdot, (a_i^S, \text{ty}_i^S), \cdot)\}_{i=1}^{ \mathcal{S} }</math> Parse <math>\mathcal{T}</math> as <math>\{(\cdot, \text{note}_i^T, C_i^T, (a_i^T, \text{ty}_i^T), \cdot)\}_{i=1}^{ \mathcal{T} }</math> for <math>\text{ty} \in \{\text{ty}_i^S\} \cup \{\text{ty}_i^T\}</math> do   <math>\Delta_{\text{ty}} \leftarrow \sum_{(a_i^S, \text{ty}) \in \mathcal{S}} a_i^S - \sum_{(a_i^T, \text{ty}) \in \mathcal{T}} a_i^T</math> return <math>(\{\text{nul}_i\}_{i=1}^{ \mathcal{S} }, \{\text{note}_i^T, C_i^T\}_{i=1}^{ \mathcal{T} }, \{\Delta_{\text{ty}}\}_{\text{ty} \in \text{Ty}})</math>  MergeTx(<math>\{\text{tx}_j\}_{j=1}^n</math>) Parse tx<math>_j</math> as <math>(\{\text{nul}_{j,i}\}_{i=1}^{ \mathcal{S} _j}, \{(\text{note}_{j,i}^T, C_{j,i}^T)\}_{i=1}^{ \mathcal{T} _j}, \{\Delta_{j,\text{ty}}\}_{\text{ty} \in \text{Ty}_j})</math> return <math>(\text{ToSet}(\bigcup_{j,i} \{\text{nul}_{j,i}\}), \text{ToSet}(\bigcup_{j,i} \{(\text{note}_{j,i}^T, C_{j,i}^T)\}),</math>   <math>\text{ToSet}(\{\sum_j \Delta_{j,\text{ty}}\}_{\text{ty} \in \bigcup_j \text{Ty}_j}))</math>  TryReceive(note, C, SK) for sk <math>\in \text{SK}</math> do   res <math>\leftarrow \text{OTA.Receive}(\text{note}, C, \text{sk})</math>   if res = <math>((a, \text{ty}), r) \neq \perp</math> then     nul <math>\leftarrow \text{OTA.NulEval}(\text{sk}, r)</math>     return (sk, nul, (a, ty), r) return <math>\perp</math> </pre>
---	--

Fig. 2. Auxiliary algorithms for Zswap. These only depend on the OTA scheme and do not use Zswap methods such as Verify or SignTx.

transactions by themselves cannot be included in the public ledger. Both have a type with a negative balance. Alice's transaction has  $\Delta_{\mathcal{E}} = -5$  and Bob's transaction has an imbalance of  $\Delta_{\mathcal{S}} = -7$ . So far Alice and Bob have not communicated. The non-malleability of the transactions allow them to publish their transactions into an exchange pool. Exchange pools may be run globally or with limited access. The first party seeing both transactions recognizes that they are complementary and is able to merge both tx $_i$  together with MergeTx and their signatures  $\sigma_i$  with MergeSig. As a prize, the merger is allowed to claim the surplus of 1\$ paid by Alice. Technically there are now 3 transactions merged. The resulting merged transaction has no imbalance and can be included in the public ledger. Then all parties get their specified outputs and store them for future transactions.

## 4.2 Zswap Correctness

**Definition 4.2** (Correctness). *A Zswap scheme is correct if OTA is correct and if for all  $\lambda \in \mathbb{N}$  and all  $p \in \text{Setup}(1^\lambda)$  it holds that:*

- (1) **Honestly generated transactions are valid:**  
For any  $\mathcal{S}, \mathcal{T}, \text{st}$  such that  $\text{CheckPTx}(\mathcal{S}, \mathcal{T}, \text{st}) = 1$ , and  $\text{tx} = \text{CompleteTx}(\mathcal{S}, \mathcal{T})$ , and for any signature  $\sigma \leftarrow \text{SignTx}(\mathcal{S}, \mathcal{T}, \text{st})$ , it holds that  $\text{Verify}(\text{st}, \text{tx}, \sigma) = 1$ .

- (2) **Honestly merged valid transactions are again valid:** For any valid st,  $n \in \mathbb{N}$  and each set of  $\{(\mathcal{S}_i, \mathcal{T}_i)\}_{i \in [n]}$  s.t. for all  $i$ ,  $\text{CheckPTx}((\mathcal{S}_i, \mathcal{T}_i), \text{st}) = 1$ , for each  $\sigma_i \in \text{SignTx}(\mathcal{S}_i, \mathcal{T}_i, \text{st})$  and  $\sigma = \text{MergeSig}(\{\sigma_i\}_{i=1}^n)$ , it holds that  $\text{Verify}(\text{st}, \text{tx}, \sigma) = 1$  where  $\text{tx} \leftarrow \text{MergeTx}(\{\text{CompleteTx}(\mathcal{S}_i, \mathcal{T}_i)\}_{i=1}^n)$ .
- (3) **Valid transactions are valid in a future state:** For any two linked states st, st' (satisfying  $\text{ev}(\text{st}, \text{st}') = 1$ ) and any transaction tx and signature  $\sigma$  with  $\text{Verify}(\text{st}, \text{tx}, \sigma) = 1$ , it holds that  $\text{Verify}(\text{st}', \text{tx}, \sigma) = 1$ , unless  $\text{st}'.\text{Nf} \setminus \text{st}.\text{Nf}$  contains any of the tx nullifiers.

The security of our scheme is based on several top-level games, plus the games we define for the OTA scheme.

## 4.3 Support Oracles

To model both security and privacy we need to introduce the oracles which will model note spending that the adversary  $\mathcal{A}$  sees, interactively. The relevant oracles,  $\mathcal{O}_{\text{KeyGen}}$ ,  $\mathcal{O}_{\text{Spend}}$  and  $\mathcal{O}_{\text{Insert}}$  are presented on Fig. 3. They are responsible for generating keys, honest transactions and inserting them into the ledger.

$\mathcal{O}_{\text{KeyGen}}$  allows generating honest public and secret keys for further use in  $\mathcal{O}_{\text{Spend}}$ .

<p><math>\mathcal{O}_{\text{Spend}}(I, O, \text{st})</math></p> <pre> <b>if</b> Spent <b>is</b> <math>\perp</math> <b>then</b> Spent <math>\leftarrow \emptyset</math> <b>assert</b> <math> I  &gt; 0 \wedge  O  &gt; 0 \wedge \exists(\text{pk}, \cdot, \cdot) \in O : \text{pk} \in \text{PK}</math> <math>\text{ptx} @ (\mathcal{S}, \mathcal{T}) \leftarrow \text{BuildPTx}(I, O, \text{st}, \text{SK})</math> <math>\sigma \leftarrow \text{SignTx}(\text{ptx}, \text{st})</math> Parse <math>\mathcal{S}</math> as <math>\{(\text{nul}_i, \cdot, \cdot, \cdot)\}_{i=1}^{ \mathcal{S} }</math>; <math>\mathcal{T}</math> as <math>\{(\cdot, \text{note}_j^{\mathcal{T}}, C_j, \cdot, \cdot)\}_{j=1}^{ \mathcal{T} }</math>; Spent := Spent <math>\cup</math> (st, I, <math>\{\text{nul}_i\}_{i=1}^{ \mathcal{S} }</math>, <math>\{(\text{note}_j^{\mathcal{T}}, C_j)\}_{j=1}^{ \mathcal{T} }</math>) <b>return</b> <math>\sigma</math> </pre> <p><math>\text{GetLog}(\text{st})</math></p> <pre> <math>\text{Ins}' \leftarrow \square</math> <b>while</b> <math>x @ (\text{st}', \text{st}, \cdot, \cdot) \in \text{Ins}</math> <b>do</b> <math>\text{Ins}' := x \parallel \text{Ins}'</math>; <math>\text{st} := \text{st}'</math> <b>assert</b> <math>\text{st} = \text{st}_0</math> <b>return</b> <math>\text{Ins}'</math> </pre>	<p><math>\mathcal{O}_{\text{Insert}}(\text{st}, \text{tx}, \sigma)</math></p> <pre> <b>if</b> Ins <b>is</b> <math>\perp</math> <b>then</b> Ins <math>\leftarrow \emptyset</math> <b>unless</b> <math>\text{st}_0 \neq \perp</math> <b>assert</b> <math>(\text{st}, \cdot, \cdot, \cdot) \in \text{Ins}</math> <b>assert</b> <math>\text{GetLog}(\text{st}) \neq \perp</math> <math>(\text{Nf} @ \{\text{nul}_i\}_{i=1}^{ \mathcal{S} }, \{\text{note}_j^{\mathcal{T}}, \cdot\}_{j=1}^{ \mathcal{T} }, \{\Delta_{\text{ty}}\}_{\text{ty} \in \text{Ty}}) \leftarrow \text{tx}</math> <b>assert</b> <math>\text{Verify}(\text{st}, \text{tx}, \sigma) = 1</math> <b>assert</b> <math>\forall \text{ty} \in \text{Ty} : \Delta_{\text{ty}} \geq 0</math> Construct <math>\text{st}'</math> s.t. <math>\text{st}'.\text{MT} \leftarrow \text{st}.\text{MT}.\text{insert}(\{\text{note}_j^{\mathcal{T}}\}_{j=1}^{ \mathcal{T} })</math> <math>\text{st}'.\text{NF} \leftarrow \text{st}.\text{NF} \cup \text{Nf}</math> <math>\text{Ins} := \text{Ins} \cup (\text{st}, \text{st}', \text{tx}, \sigma)</math> <b>if</b> <math>\text{st}_0 = \perp</math> <b>then</b> <math>\text{st}_0 \leftarrow \text{st}</math> </pre> <p><math>\mathcal{O}_{\text{KeyGen}}()</math></p> <pre> <b>if</b> SK or PK <b>is</b> <math>\perp</math> <b>then</b> SK, PK <math>\leftarrow \emptyset</math> (sk, pk) <math>\leftarrow \text{KeyGen}()</math> SK := SK <math>\parallel</math> sk; PK := PK <math>\parallel</math> pk <b>return</b> pk </pre>
---	---

Fig. 3. Oracles for the security and privacy experiments.

$\mathcal{O}_{\text{Spend}}$  models leakage of honestly generated transactions, including unbalanced “offer” transactions. Each  $\mathcal{O}_{\text{Spend}}$  query allows spending some notes, and successful spend logs are recorded in Spent. An adversary can specify any possible state it likes as long as the request is valid with respect to this state. The requirement that  $O$  has at least one honest output is a modelling artefact, and is explained further in the “anti-theft” section.

$\mathcal{O}_{\text{Insert}}$  models “recording” balanced transaction in the ledger. Branching is possible inside  $\mathcal{O}_{\text{Insert}}$  —  $\mathcal{A}$  can append an honest transaction to any of the recorded states, but still one can only spend notes through  $\mathcal{O}_{\text{Insert}}$  if their nullifier has not been used in the same branch. Moreover, the  $\text{st}_0$  variable is shared between the oracle records as the first “root” state that  $\mathcal{A}$  can initialize the oracles with.  $\mathcal{O}_{\text{Insert}}$  only accepts states that are eventually linked with this root  $\text{st}_0$ .

#### 4.4 Anti-Theft and Non-Malleability

The notion of anti-theft we describe here is our main non-malleability notion. Intuitively it says that an adversary  $\mathcal{A}$  can only merge transactions, but cannot break honest transactions apart, or modify them in any other way. More concretely, for any honest tx that  $\mathcal{A}$  sees,  $\mathcal{A}$  cannot submit  $\text{tx}^*$  which contains any of the (1) input nullifiers of tx; or (2) honest output notes of tx; without merging the whole tx into  $\text{tx}^*$ .

The anti-theft game is modelled by allowing  $\mathcal{A}$  to interact with the oracles described before:  $\mathcal{A}$  can generate new keys, produce unbalanced honest transactions (spend), and submit (insert) transactions (which we want to prove to be only honest, dishonest, or merges of the two types). The adversary then returns a state  $\text{st}^*$  which is examined. The challenger searches through

#### Anti-Theft $_{\mathcal{A}}(1^\lambda)$

```

 $p \leftarrow \text{Setup}(1^\lambda)$ 
 $\text{st}^* \leftarrow \mathcal{A}^{\mathcal{O}_{\text{KeyGen}}, \mathcal{O}_{\text{Spend}}, \mathcal{O}_{\text{Insert}}}(p)$ 
for  $(\text{st}, \cdot, \text{tx}, \cdot) \in \text{GetLog}(\text{st}^*)$  do
 $(\cdot, (\text{Nf}^{\mathcal{A}}, M^{\mathcal{A}})) \leftarrow \text{SplitTx}(\text{st}, \text{tx})$ 
 $M^{\mathcal{A}'} \leftarrow \{(\text{note}, C) \in M^{\mathcal{A}} \mid \text{TryReceive}(\text{note}, C, \text{SK}) \neq \perp\}$ 
if  $\exists(\cdot, \cdot, \text{Nf}, M) \in \text{Spent} : M^{\mathcal{A}'} \cap M \neq \emptyset \vee \text{Nf}^{\mathcal{A}} \cap \text{Nf} \neq \emptyset$ 
then
return 1
return 0

```

$\text{SplitTx}(\text{st}, \text{tx})$

```

 $(\text{Nf} @ \{\text{nul}_i\}_{i=1}^{|\mathcal{S}|}, M @ \{\text{note}_j^{\mathcal{T}}, \cdot\}_{j=1}^{|\mathcal{T}|}, C) \leftarrow \text{tx}$ 
 $\text{Nf}_{\mathcal{H}}, M_{\mathcal{H}} \leftarrow \emptyset$ ;  $\text{Nf}_0 \leftarrow \text{Nf}$ 
for  $(\text{st}', \cdot, \text{Nf}' @ \{\text{nul}'_i\}_{i=1}^{|\mathcal{S}'|}, M' @ \{\text{note}'_j, C'_j\}_{j=1}^{|\mathcal{T}'|}) \in \text{Spent}$  do
if  $\text{st}'.\text{MT} = \text{st}.\text{MT} \wedge M' \subset M \wedge \text{Nf}' \subset \text{Nf}_0$  then
 $\text{Nf}_{\mathcal{H}} := \text{Nf}_{\mathcal{H}} \cup \text{Nf}'$ ;  $M_{\mathcal{H}} := M_{\mathcal{H}} \cup M'$ 
 $\text{Nf}_0 := \text{Nf}_0 \setminus \text{Nf}'$ 
return  $(\text{Nf}_{\mathcal{H}}, M_{\mathcal{H}}), (\text{Nf} \setminus \text{Nf}_{\mathcal{H}}, M \setminus M_{\mathcal{H}})$ 

```

Fig. 4. The anti-theft experiment.

the log of inserted transactions. The adversary wins if a transaction in the log contains an incomplete part of (an honest) transaction earlier returned by  $\mathcal{O}_{\text{Spend}}$  to  $\mathcal{A}$ , but *used* only partially. To do that, it calls a sub-function SplitTx which locates complete “honest subtransactions” returned by  $\mathcal{O}_{\text{Spend}}$  in tx, and returns these honest subtransactions as first argument, and the remaining parts of  $\text{tx}^*$ . If the challenger recognizes nullifiers or output notes in this remaining part, that were created in  $\mathcal{O}_{\text{Spend}}$ , it means that  $\mathcal{A}$  managed to deconstruct it, changing the output of  $\mathcal{O}_{\text{Spend}}$ , or stealing an honest nullifier.

**Definition 4.3** (Anti-Theft). *A Zswap scheme is protecting against theft, if for any  $\lambda \in \mathbb{N}$  and any PPT adversary  $\mathcal{A}$ ,  $\Pr[\text{Anti-Theft}_{\mathcal{A}}(1^\lambda) = 1] \leq \text{negl}(\lambda)$  where the Anti-Theft $_{\mathcal{A}}(1^\lambda)$  game is defined in Fig. 4.*

Essentially, anti-theft captures non-malleability of output notes and their ciphertexts, which is in practice guaranteed by NIZK SE (that implies instance binding). If  $\mathcal{A}$  can maul an output (e.g. change its value, or destination) and still submit the tx to  $\mathcal{O}_{\text{Insert}}$  successfully, it wins the game, since `SplitTx` will not locate tx as being in `Spent`, and thus the game challenger  $\mathcal{C}$  will catch  $\mathcal{A}$  spending a nullifier that was in `Spent` but not located by `SplitTx`.

The requirement that  $O$  in  $\mathcal{O}_{\text{Spend}}$  contains at least one honest output is needed so that `SplitTx` can uniquely identify honest sub-transactions<sup>7</sup>. In other words, note uniqueness is necessary for the anti-theft game to make sense. This does not limit  $\mathcal{A}$  from creating unbalanced input-only transactions, since  $\mathcal{A}$  can still request to include a single zero-valued output.

The set  $\text{Nf}_0$  in `SplitTx` is needed since without it honest sub-transactions can be counted twice. E.g. let  $\text{tx}_0, \text{tx}_1$  be two honest transactions with the same input nullifiers  $\text{Nf}$ , but completely different outputs  $M_0, M_1$  (each containing at least one honest output note). Without  $\text{Nf}_0$ , for tx containing  $\text{Nf}$  and  $M_0 \cup M_1$ , `SplitTx` will detect both  $\text{tx}_0, \text{tx}_1$ , and the winning condition will not be triggered. We, on the other hand, want anti-theft to prevent this: with  $\text{Nf}_0$  one of  $\text{tx}_i$  will be considered honest, and the honest output note in  $\text{tx}_{1-i}$  will trigger the anti-theft winning condition. We do not need to similarly count  $M$  since output notes are unique — adding  $M = M \setminus M'$  into `SplitTx` simply does not change the behaviour of the game.

Finally, it is critical to filter  $M^{A'}$  from  $M^A$ . Otherwise,  $\mathcal{A}$  would be able to trivially win the game by triggering  $M^A \cap M \neq \emptyset$  in the following way:

1.  $\mathcal{A}$ , through  $\mathcal{O}_{\text{Spend}}$ , requests an honest spend to an adversarial public key  $\text{pk}^A$ ;
2.  $\mathcal{A}$  obtains the proof and corresponding  $\text{note}^*$ , receives it, and creates a completely different, adversarial transaction  $\text{tx}^*$  on its own, where  $\text{note}^*$  is an output;
3. the game will not find any honest subtransaction in  $\text{tx}^*$  and thus  $\text{note}^* \in M \cap M^A$ .

<sup>7</sup> Without this requirement, and without adding extra marker information to notes, identifying honest sub-transactions within `SplitTx` would take exponential time.

### Balance $_{\mathcal{A}, \mathcal{E}_A}(1^\lambda)$

```

p ← Setup(1λ)
(st, I0, I, SK*) ← AKeyGen · OSpend · OInsert(p)
assert {note0,i} ∈ I0 are distinct, and st0.MT contains only
these notes.
assert {notei} ∈ I are distinct, in st.MT
for (notei, Ci) ∈ I0 do
  (·, nul, (a, ty), ·) ← TryReceive(notei, Ci, SK*)
  assert nul ∉ st0.NF
  v0[ty] := v0[ty] + a
for (notei, Ci) ∈ I do
  (·, nul, (a, ty), ·) ← TryReceive(notei, Ci, SK*)
  assert nul ∉ st.NF
  vA[ty] := vA[ty] + a
vH-, vH+ ← (ty ↦ 0) % map with default value 0
Ins' ← GetLog(st) % Also makes sure ev(st0, st)
for all (st, ·, tx, σ) ∈ Ins' do
  ({nuli}i=1|S|, {notejT}j=1|T|, {Δty}ty ∈ Ty) ← tx
  (txH, ·) ← SplitTx(st, tx)
  for (Nf, M) ∈ txH do
    Find(st, I, Nf, M) ∈ Spent
    for (note', C') ∈ I do
      (·, (a, ty), ·) ← TryReceive(note', C', SK)
      vH-[ty] := vH-[ty] + a
  for (·, M) ∈ tx do
    for (note', C') ∈ M do
      res ← TryReceive(note', C', SK)
      if res = (·, (a, ty), ·) then vH+[ty] := vH+[ty] + a
if ∃ty ∈ T : vA[ty] > v0[ty] + vH-[ty] - vH+[ty] then
  return 1
else return 0

```

Fig. 5. The balance experiment

## 4.5 Balance

The balance property presented next captures that transactions distribute underlying coins properly, that is, the adversarial balance per type only changes “predictably”, by the adversary receiving or sending these coins. In particular, the balance game forbids malicious conversion between asset types, coin forging, and double spending.

Formally, the property is modelled as a game (Fig. 5) of  $\mathcal{A}$  “against” the honest parties, in which  $\mathcal{A}$  has to prove that it has more coins that it could have obtained honestly. It goes as follows. The challenger allows  $\mathcal{A}$  to interact with the oracles, and asks  $\mathcal{A}$  to present two sets of values:  $I_0$  and  $I$ , together with a set of secret keys  $\text{SK}^*$ . Unlike in other games,  $\text{st}_0$  in balance must contain only adversarial coins, and  $I_0$  must be all the unspent corresponding notes (receivable with  $\text{SK}^*$ ). The other set  $I$  is the set of unspent adversarial coins from the new state  $\text{st}$  (again receivable using  $\text{SK}^*$ ) that  $\mathcal{A}$  claims to have an illicitly produced coins, breaking the balance condition. This balance condition is computed next in the following manner. First,  $v_0$  is set to the sum of all values in  $I_0$ , and  $v_A$  as a sum of all values in  $I$ . Then, by traversing the history of transactions from  $\text{st}_0$  to  $\text{st}$ , the game computes the following two values:

1.  $v_{\mathcal{H}_-}$  is the sum of all honest inputs in transactions; and
2.  $v_{\mathcal{H}_+}$  is the total coins received by honest parties.

The final condition checks that  $\mathcal{A}$  cannot show in  $\mathbf{st}$  more coins than:

1. it had in  $\mathbf{st}_0$ , and
2. that were sent by honest parties, minus received by honest parties.

This last difference may be both positive and negative, but importantly, the balance of  $\mathcal{A}$ , per type, is “tied” to the balance of honest parties, no extra coins can be produced.

**Definition 4.4** (Balance). *A Zswap scheme is balanced if for all PPT adversaries  $\mathcal{A}$  there exists a PPT extractor  $\mathcal{E}_{\mathcal{A}}$  such that  $\Pr[\text{Balance}_{\mathcal{A}, \mathcal{E}_{\mathcal{A}}}(1^\lambda) = 1] \leq \text{negl}(\lambda)$  with the game defined in Fig. 5.*

Note that this balance property does not prevent an external mechanism to change the state. This is useful for smart contract based minting of new types.

## 4.6 Privacy

The privacy game captures secrecy of coin transfers by means of an indistinguishability experiment. To model the game we will first introduce the notion of a transaction instruction tree  $T$ . Such a variable width tree has as its leaf  $i$  either

1. the instructions to construct an honest transaction  $(\{\text{note}_{i,j,C_{i,j}}\}_{j=1}^{|S|_i}, \{(\text{pk}_{i,j}, a_{i,j}, \text{ty}_{i,j})\}_{j=1}^{|T|_i})$ , similarly to the input to  $\mathcal{O}_{\text{Spend}}$ , or
2. a fully adversarial transaction  $(\text{tx}_i, \sigma_i)$ .

Its intermediate leaves are empty, and merely represent how children transactions must be merged.

We will also need to decide what a transaction resulting from a merge according to  $T$  leaks. We formalize this notion by defining the tree equivalency relation, formally on Fig. 6 as follows:  $\text{EquivTree}(T_0, T_1) = 1$  if

1. The imbalance of amounts in each type is equal,  $B_1$ .
2. The number of input notes is equal (published nullifiers), and the number of output notes is equal. For honest  $B_2^S \wedge B_2^T$  leaves, as it implicitly holds for adversarial leaves due to the next property  $B_3$ .
3. Malicious offers need to be the same in both trees, but maybe not at same positions,  $B_3$ .
4. For all the honest nullifiers  $\text{Nf}$  that  $\mathcal{A}$  receives via  $\mathcal{O}_{\text{Spend}}$ , if any of the related notes are included in the honest leaves of a tree, these notes must be included in both trees,  $B_4$ .

5. The adversarial output instructions of honest leaves are the same in both trees,  $B_5$ .

The privacy notion itself asks  $\mathcal{A}$  to present two equivalent trees, and then builds a single transaction for each tree which both must not fail. It returns one merged transaction. The adversary wins the game if it can decide which tree was used. To illustrate the notion with a single example, imagine the tree to contain a single leaf node in both cases, in the first case spending a single note with  $X$  coins, sending 1 to Alice and  $N - 1$  to Bob, and in the second case spending a single note of  $Y$  coins, sending  $Y - 2$  to Alice and 2 to Bob. Such trees are equivalent according to our definition, and thus  $\mathcal{A}$  should not be able to decide which transaction it sees.

**Definition 4.5** (Transaction Privacy). *A Zswap scheme has private transactions, if for all PPT adversaries  $\mathcal{A}$  it holds that:*

$$|\Pr[\text{Privacy}_{\mathcal{A}}^0(1^\lambda) = 1] - \Pr[\text{Privacy}_{\mathcal{A}}^1(1^\lambda) = 1]| \leq \text{negl}(\lambda)$$

with  $\text{Privacy}_{\mathcal{A}}^b(1^\lambda)$  defined in Fig. 6.

## 5 Zswap Protocol

We present the Zswap protocol in Fig. 7. It extends the OTA scheme (constructed in Appendix B) and additionally utilizes a sparse homomorphic commitment scheme (sparse Pedersen):  $\text{SHC} = (\text{ComSetup}, \text{Commit})$ . and Arguments of Knowledge:  $\text{AoK} = (\text{AoK.Setup}, \text{AoK.Prove}, \text{AoK.Verify}, \text{AoK.Sim})$ .

The high-level idea is to create a transaction which has separate inputs and outputs handled by the OTA scheme. A transaction links them together through SHC commitments. Each input and output has a corresponding SHC commitment with an equal amount and type. This equivalency is assured by an AoK for each input and each output. The output AoK additionally assure non-malleability for the note ciphertext. The input AoKs enforce that the transaction creator possessed the secret key to authorize the spending and prove that the published nullifier is correct. This is captured by two NIZK languages.

The first one authenticates a valid spend — it says that the (rerandomized) SHC commitment  $\text{com}^S$  is well-formed and contains the same value and type as a note

Privacy $_{\mathcal{A}}^b(1^\lambda)$	EquivTree( $T_0, T_1, \text{st}$ )
<pre> p ← Setup(1<sup>λ</sup>) O = {O<sub>KeyGen</sub>, O<sub>Spend</sub>, O<sub>Insert</sub>} (st, T<sub>0</sub>, T<sub>1</sub>) ← A<sup>O</sup>(p) assert GetLog(st) ≠ ⊥ tx<sub>0</sub> ← EvalTree(st, T<sub>0</sub>) tx<sub>1</sub> ← EvalTree(st, T<sub>1</sub>) assert tx<sub>0</sub> ≠ ⊥ ∧ tx<sub>1</sub> ≠ ⊥ assert EquivTree(T<sub>0</sub>, T<sub>1</sub>, st) = 1 b' ← A<sup>O</sup>(tx<sub>b</sub>) return b' </pre>	<pre> Parse each leaf<sub>b,i</sub> of each T<sub>b</sub> as either an adversarial leaf leaf<sup>A</sup> (tx<sub>b,i</sub> ⊗ ({nul<sub>b,i,j</sub>}<sup> S <sub>b,i</sub></sup>, {(note<sub>b,i,j</sub><sup>T</sup>, C<sub>b,i,j</sub><sup>T</sup>)<sup> T <sub>b,i</sub></sup>}, {Δ<sub>b,i,ty</sub>}<sup>ty ∈ Ty<sub>b,i</sub></sup>), σ<sub>b,i</sub>) ← leaf<sup>A</sup><sub>b,i</sub> or an honest leaf leaf<sup>H</sup> ({(note<sub>b,i,j</sub><sup>S</sup>, C<sub>b,i,j</sub><sup>S</sup>)<sup> S <sub>b,i</sub></sup>, {(pk<sub>b,i,j</sub><sup>T</sup>, a<sub>b,i,j</sub><sup>T</sup>, ty<sub>b,i,j</sub><sup>T</sup>)<sup> T <sub>b,i</sub></sup>}) ← leaf<sup>H</sup><sub>b,i</sub> for b ∈ {0, 1}, leaf<sup>H</sup><sub>i</sub> ∈ T<sub>b</sub> do for j ∈ [ T<sub>b,i</sub> ] do (nul<sub>b,i,j</sub>, (a<sub>b,i,j</sub><sup>S</sup>, ty<sub>b,i,j</sub><sup>S</sup>), r<sub>b,i,j</sub>) ← TryReceive(note<sub>b,i,j</sub><sup>S</sup>, C<sub>b,i,j</sub><sup>S</sup>, SK) for ty ∈ {ty<sub>b,i,j</sub><sup>S</sup>}<sup> T <sub>b,i</sub></sup> do Δ<sub>b,i,ty</sub> ← ∑<sub>j:ty<sub>b,i,j</sub><sup>S</sup>=ty</sub> a<sub>b,i,j</sub><sup>S</sup> - ∑<sub>j:ty<sub>0,i,j</sub><sup>S</sup>=ty</sub> a<sub>0,i,j</sub><sup>T</sup> B<sub>1</sub> ← ∀ty : ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>0</sub></sub> Δ<sub>0,i,ty</sub> = ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>1</sub></sub> Δ<sub>1,i,ty</sub> B<sub>2</sub><sup>S</sup> ← ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>0</sub></sub>  S <sub>0,i</sub> = ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>1</sub></sub>  S <sub>1,i</sub> B<sub>2</sub><sup>T</sup> ← ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>0</sub></sub>  T <sub>0,i</sub> = ∑<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>1</sub></sub>  T <sub>1,i</sub> B<sub>3</sub> ← ∪<sub>leaf<sup>A</sup><sub>i</sub> ∈ T<sub>0</sub></sub> {(tx<sub>0,i</sub>, σ<sub>0,i</sub>)} = ∪<sub>leaf<sup>A</sup><sub>i</sub> ∈ T<sub>1</sub></sub> {(tx<sub>1,i</sub>, σ<sub>1,i</sub>)} Nf ← ∪<sub>(·, ·, Nf<sub>i</sub>, ·) ∈ Spent</sub> Nf<sub>i</sub> B<sub>4</sub> ← Nf ∩ (∪<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>0</sub></sub> nul<sub>0,i,j</sub>) = Nf ∩ (∪<sub>leaf<sup>H</sup><sub>i</sub> ∈ T<sub>1</sub></sub> nul<sub>1,i,j</sub>) for b ∈ {0, 1}, leaf<sup>H</sup><sub>i</sub> ∈ T<sub>b</sub> do O<sub>b</sub><sup>A</sup> ← ∪<sub>pk<sub>b,i,j</sub><sup>T</sup> ∈ PK</sub> (pk<sub>b,i,j</sub><sup>T</sup>, a<sub>b,i,j</sub><sup>S</sup>, ty<sub>b,i,j</sub><sup>S</sup>) B<sub>5</sub> ← O<sub>0</sub><sup>A</sup> = O<sub>1</sub><sup>A</sup> return B<sub>1</sub> ∧ B<sub>2</sub><sup>S</sup> ∧ B<sub>2</sub><sup>T</sup> ∧ B<sub>3</sub> ∧ B<sub>4</sub> ∧ B<sub>5</sub> </pre>
<pre> EvalTree(st, T) for all leaf<sub>i</sub> ∈ T do if leaf<sub>i</sub> = (tx, σ) then % Validate adversarial txs assert Verify(st, tx, σ) = 1 % Replace instruction leaves with real txs if leaf<sub>i</sub> = (I, O) then ptx ← BuildPTx(I, O, st, SK) tx<sub>i</sub> ← CompleteTx(ptx) σ<sub>i</sub> ← SignTx(ptx, st) Replace leaf<sub>i</sub> by (tx<sub>i</sub>, σ<sub>i</sub>) in T % Fold T into a single root node by merging while ∃ node N in T with only {(tx<sub>i</sub>, σ<sub>i</sub>)<sup>c</sup><sub>i=1</sub>} as children do Replace node with merged transactions of its children: N ← (MergeTx({tx<sub>i</sub>}<sup>c</sup><sub>i=1</sub>, MergeSig({σ<sub>i</sub>}<sup>c</sup><sub>i=1</sub>))) Remove its children return T </pre>	

Fig. 6. Transaction privacy experiment

in the Merkle tree of state  $\text{st}$  with the given nullifier  $\text{nul}$ .

$$\mathcal{L}^{\text{spend}} = \left\{ (\text{st}, \text{nul}, \text{com}^S) \mid \exists (\text{path}, \text{sk}^S, a^S, \text{ty}^S, r^S, \text{rc}^S) : \right. \\
(\text{st.MT}[\text{path}], \text{nul}; \text{sk}^S, (a^S, \text{ty}^S), r^S) \in \mathcal{L}^{\text{nul}} \wedge \\
(\text{st.MT}[\text{path}]; \text{OTA.P}(\text{sk}^S), (a^S, \text{ty}^S), r^S) \in \mathcal{L}^{\text{open}} \wedge \\
\left. \text{com}^S = \text{Commit}(\text{ty}^S, a^S; \text{rc}^S) \right\}$$

The second language  $\mathcal{L}^{\text{output}}$  is even simpler. It claims that the two output commitments, the real (which is contained inside the output note) and the randomized one, contain the same value of the same type. To prevent overflows in the homomorphic commitments, we include a range proof where  $\beta$  is chosen small enough in relation to the group order (maximum inputs and outputs times  $2^\beta < |\mathbb{G}|$ ). Thereby we assume integer amounts in subsequent arguments.

$$\mathcal{L}^{\text{output}} = \left\{ (\text{note}^T, C^T, \text{com}^T) \mid \exists (\text{pk}^T, a^T, \text{ty}^T, r^T, \text{rc}^T) : \right. \\
(\text{st.MT}[\text{path}]; \text{pk}, (a^T, \text{ty}^T), r^T) \in \mathcal{L}^{\text{open}} \wedge \\
\text{com}^T = \text{Commit}(\text{ty}^T, a^T; \text{rc}^T) \\
\left. a^T \in \{0, \dots, 2^\beta - 1\} \right\}$$

Note that the ciphertext  $C^T$  is not referred to in the relation, but when used with a simulation extractable (SE) NIZK, realizes a Signature of Knowledge. I.e. every proof is bound to a specific ciphertext and is invalid for any other ciphertext. The transaction signature  $\sigma$

then contains a proof  $\pi_i^S$  for each input together with the SHC  $\text{com}_i^S$  and for each output a signature  $\pi_i^T$  and the SHC  $\text{com}_i^T$ . The last component of the signature is the aggregated randomness of the commitments which, together with the  $\Delta_{\text{ty}}$  imbalance, allows verification.

For a full transaction verification, all proofs and signatures contained in  $\sigma$  must be valid and the published nullifiers must be unique regarding the set of nullifiers in state  $\text{st}$ .

With the transaction signature  $\sigma$  having a separate proof for each input and output, it is possible to merge transactions by calculating the union of their proofs. To maintain the verifiability of the commitments, the randomness of the merged transactions is added. The irreversible addition operation then prevents future parties to unmerge a transaction if they have not seen the separate parts beforehand. To maintain the anonymity, we order inputs and outputs canonically after each merge.

As a remark, the aggregated randomness in a transaction may be replaced by a proof of knowledge. Like the binding signature of Sapling, this finalizes a transaction such that it can no longer be merged with others.

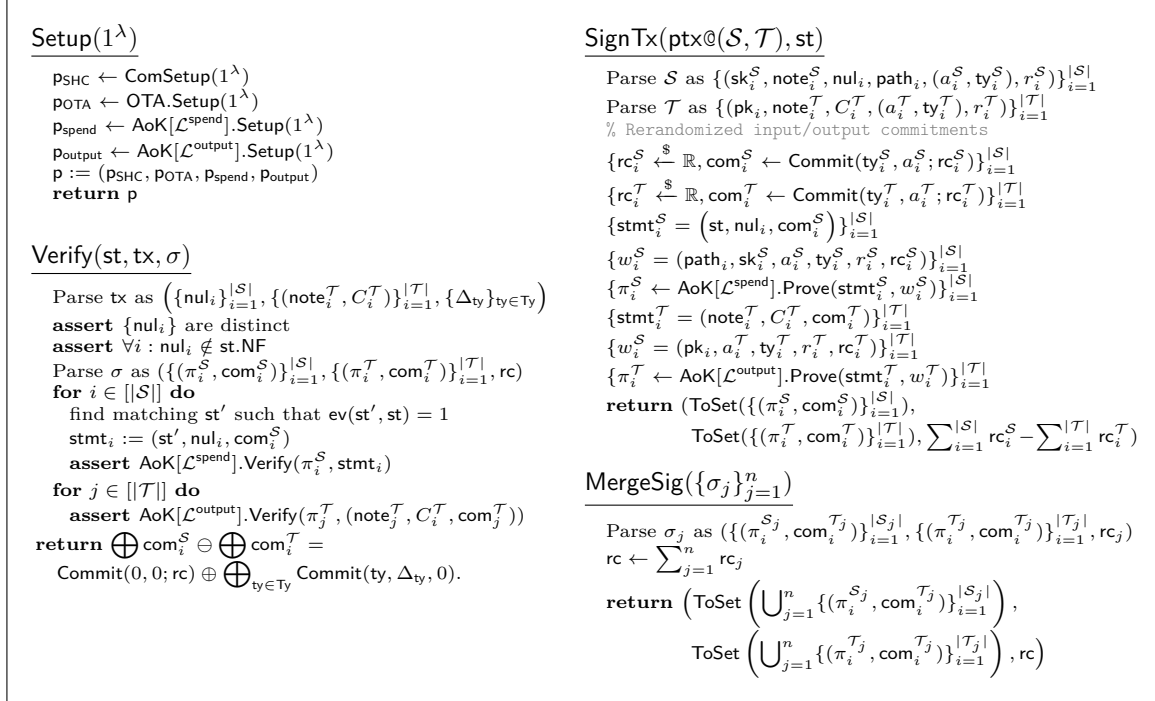


Fig. 7. The Zswap Construction

## 6 Security Proof

In this section we prove the main three security properties of Zswap construction on Fig. 7 we introduced in Section 4.

**Theorem 6.1** (Anti-Theft). *The Zswap protocol prevents theft (Definition 4.3), assuming OTA security, NIZK zero-knowledge, NIZK simulation-extractability, and SHC binding and HID-OR.*

*Proof (Sketch).* When  $\mathcal{A}$  triggers the winning condition of the anti-theft game with an adversarial transaction  $\text{tx}^*$ , there exists a note or a nullifier taken from some honest  $\mathcal{O}_{\text{Spend}}$  query  $\hat{E}$ , such that not all nullifiers and notes from  $\text{tx}$  were included in  $\text{tx}^*$ . By NIZK simulation-extractability, the proofs that were produced in  $\hat{E}$  “bind” together notes and nullifiers with the corresponding input and output commitments. This means  $\mathcal{A}$  uses some commitments from  $\text{tx}$ , but avoids some other. Given commitment binding, the values we extract for commitments of  $\text{tx}^*$  (we can extract them from NIZKs) are the same as the values committed in  $\text{tx}$ .

From this point we can build a reduction  $\mathcal{B}$  to HID-OR.  $\mathcal{B}$  guesses a commitment  $C_X$  that is present in both  $\text{tx}, \text{tx}^*$ , and  $C_Y$  that is only present in  $\text{tx}$ . It asks the HID-OR challenger  $\mathcal{C}$  for either

1. two commitments corresponding to the values  $(a_1, \text{ty}_1), (a_2, \text{ty}_2)$  as they should be honestly; or
2. the values swapped as  $(a_2, \text{ty}_1), (a_1, \text{ty}_2)$  if both indices correspond to inputs or both to outputs, and
3. swapped with negation  $(-a_2, \text{ty}_2), (-a_1, \text{ty}_1)$  otherwise.

This way of embedding guarantees the HID-OR requirement that  $\Delta_{\text{ty},0} = \Delta_{\text{ty},1}$ . Then  $\mathcal{B}$  simulates the two NIZKs corresponding to  $C_X, C_Y$ . When  $\mathcal{A}$  presents  $\text{tx}^*$ ,  $\mathcal{B}$  will extract the randomness  $\text{rc}_i$  for all commitments except for  $C_X$ , and thus because  $\text{tx}^*$  also includes the joint randomness  $\text{rc}^*$ ,  $\mathcal{B}$  can compute the randomness  $\text{rc}_X$  for  $C_X$  (as  $\text{rc}^* - \sum \text{rc}_i$ ). Given  $\text{rc}_X$ ,  $\mathcal{B}$  can check whether  $C_X$  contains  $(a_1, \text{ty}_1)$  or the “swapped” value, and thus wins HID-OR.  $\square$

We present the detailed proof of anti-theft in the full version of this paper.

**Theorem 6.2** (Privacy). *The Zswap protocol is private (Definition 4.5), if NIZK is zero-knowledge, Pedersen commitments are hiding, and OTA is private and satisfies nullifier pseudorandomness.*

*Proof.* First, we note that changing the order of merge in the tree  $T$  does not affect the resulting merge transaction. Let  $T'$  be a variant of  $T$  where any two transactions are swapped — then  $\text{EquivTree}(T, T') = 1$  and

$\text{EvalTree}(\text{st}, T) = \text{EvalTree}(\text{st}, T')$ . The first statement can be verified by manually checking all the predicates and making sure they are indifferent to the order of the leaves. As for the evaluation equality statement, first note that all the leaves will be processed in the same way and with respect to the same  $\text{st}$  irrespectively of their order. So we only need to argue that  $\text{MergeTx}$  and  $\text{MergeSig}$  are commutative — which is trivial since they are only uniting sets and taking sums (which are commutative operations on their own).

Hence we can represent  $T$  as a merge of two transactions, coming from subtrees  $T_{\mathcal{A}}$  and  $T_{\mathcal{H}}$ .  $T_{\mathcal{A}}$  contains exactly the same leaves for both  $T_0$  and  $T_1$  if they are equivalent. Therefore, we only need to show the indistinguishability of transactions  $\text{tx}_{\mathcal{H},0}$  and  $\text{tx}_{\mathcal{H},1}$ , created from  $T_{\mathcal{H},0}, T_{\mathcal{H},1}$  correspondingly. If  $T_{\mathcal{H}}$  is empty in one case (contains no leaves), it must be empty in the other case, due to the restriction on input and output size,  $B_2$ , in  $\text{EquivTree}$ . And thus if  $\text{tx}_{\mathcal{H},b}$  is empty,  $\text{tx}_{\mathcal{H},1-b}$  should be empty too for each  $b \in \{0, 1\}$ , in which case the privacy proof is trivial, since adversarial transactions are constructed in exactly the same manner in both worlds. Therefore, assume that there is at least one input or output in an honest transaction in both cases. Next, observe that:

$$\begin{aligned} \text{MergeSig}(\{\text{SignTx}(\text{BuildPTx}(I_i, O_i, \text{st}, \text{SK}), \text{st})\}_i) = \\ \text{SignTx}(\text{BuildPTx}(I' = \bigcup I_i, O' = \bigcup O_i, \text{st}, \text{SK}), \text{st}) \end{aligned}$$

Because of this homomorphic property (and a similar one for  $\text{MergeTx}$ ), we can assume that both  $\text{tx}_{\mathcal{H},b}$  and their signatures are just a direct output of  $\text{CompleteTx}$  and  $\text{SignTx}$  (w.r.t.  $(I', O')$ ), and no merging is involved. The form of  $\text{tx}_{\mathcal{H},b}$  that  $\mathcal{A}$  receives is

$$\begin{aligned} (\{\text{nul}_i\}^{|\mathcal{S}|}, \{\text{note}_i, C_i\}^{|\mathcal{T}|}, \{\Delta_{\text{ty}}\}, \\ \sigma_i @ (\{\{\pi_i^{\mathcal{S}}, \text{com}_i^{\mathcal{S}}\}^{|\mathcal{S}|}, \{\{\pi_i^{\mathcal{T}}, \text{com}_i^{\mathcal{T}}\}^{|\mathcal{T}|}, \text{rc}) \end{aligned}$$

We first argue informally why this transaction looks the same for both trees.

1. Set sizes: number of inputs and outputs  $|\mathcal{S}|$  and  $|\mathcal{T}|$  are the same, guaranteed by  $B_2$ , and by the fact these dimensions are just summed when transactions are merged.
2. Size and content of the  $\{\Delta_{\text{ty}}\}$  set is the same by  $B_1$ .
3. The set of (honest) input nullifiers  $\{\text{nul}_i\}^{|\mathcal{S}|}$  contains: nullifiers that  $\mathcal{A}$  received from  $\mathcal{O}_{\text{Spend}}$  — these are the same by  $B_4$ ; and nullifiers unknown previously to  $\mathcal{A}$  — deterministic but indistinguishable by nullifier pseudorandomness.
4.  $\{\{\text{note}_i, C_i\}^{|\mathcal{T}|}\}$  are either: belonging to adversarial keys in which case both the notes and ciphertexts

have the same distribution, as guaranteed by  $B_5$  — these notes contain the same values, but created in both trees with uniformly sampled randomness, so are indistinguishable; the same applies to the adversarial ciphertexts. Or the notes belong to the honest keys, in which case they are (together with the ciphertexts) indistinguishable by the OTA privacy.

5. In the signatures  $\sigma_i$ :

- (5.1) By zero-knowledge, proofs  $\pi_i^{\mathcal{S}}$  and  $\pi_i^{\mathcal{T}}$  can be simulated in both worlds, so they are indistinguishable.
- (5.2) Intermediary commitments  $\text{com}_i^{\mathcal{S}}, \text{com}_i^{\mathcal{T}}$  can contain just zero (except for one chosen commitment for each type which must contain  $\Delta_{\text{ty}}$  to balance out the public imbalance), and this is indistinguishable by  $\text{HID-OR}$  — commitment hiding with open randomness (by a variation with  $n$  elements involved simultaneously, as described in Lemma 2.1).

The formal reduction is described in Appendix D.  $\square$

**Theorem 6.3** (Balance). *The Zswap protocol satisfies the Balance property (Definition 4.4) by anti-theft, NIZK ZK, NIZK SE, and OTA Security.*

*Proof (Sketch).* Intuitively, the property reduces to:

1. NIZK SE (implying knowledge soundness (KS)) and commitment homomorphism (every valid transaction is balanced);
2. Anti-theft, since  $\mathcal{A}$  can use only its own notes; and binding, since  $\mathcal{A}$  can only open output commitments in a single way when spending them.

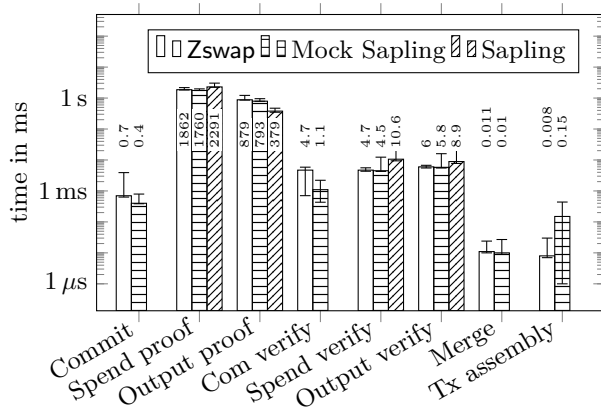
The first guarantees that transactions do not break per-type balance and do not introduce any coins out of thin air etc. The second proves that in the long-term  $\mathcal{A}$  can only move funds through such honest transactions. The full proof is presented in Appendix E.  $\square$

## 7 Implementation

To show that our construction is practical, we developed a prototype implementation available online<sup>8</sup>. We use the rust framework ark-works<sup>9</sup> and their Groth16 SNARK library. For an efficient hashing circuit, we use Poseidon [15]. As the construction is similar to the Zcash Sapling version, we use this as performance com-

<sup>8</sup> <https://anonymous.4open.science/r/zswap-E331>

<sup>9</sup> <https://arkworks.rs>



**Fig. 8.** Comparison of our protocol to our Mock Sapling implementation with the same hash function and the original Zcash.

parison. Our implementation differs in various aspects that aren't material to the core protocol, including the SNARK implementation, hash functions used, and commitments used in the coin commitment Merkle tree. The production Zcash Sapling implementation<sup>10</sup> uses less performant variants for backwards compatibility and integration purposes. A direct performance comparison to Zcash is therefore misleading, and we introduce an artificial "Mock Sapling" code base, that re-implements the Zcash Sapling cryptographic protocol but uses the same primitives as we do. We achieve this by removing the types and type commitments from our construction. The comparison then allows us to detect the performance impact of our changes to the cryptographic protocol compared to Sapling without the measurement error caused by different implementations.

For spend proofs, our protocol requires 25,416 gates, 10% more than Mock Sapling (23,084). For output proofs, we use 14,852 gates, 18% more (12,520). In both cases the difference of 2,332 gates consists of a hash-to-curve operation to associate the type with a curve point. The rest of the constraints further break down into two dominant groups: Symmetric cryptography (hashes, commitments, and Merkle tree verification), and group operations related to the binding signature. For spend proofs, symmetric operations dominate with 66% of constraints (16,935), primarily from the Merkle tree verification (using a tree of height 32). The group operations cover almost all of the remaining 24% of constraints (with 6,114). For output proofs, group operations make up 41% of the constraints (with the same 6,114 constraints), and symmetric operations 43% of constraints (6,403).

The median, min and max runtimes of 30 executions on a 6<sup>th</sup> generation i7 CPU@2.7GHz are summarized in Figure 8. The times are for creating Commitments, Spend proofs and Output proofs. For verification, we measure the homomorphic commitment comparison, the Spend proof verification and the Output proof verification. Without comparison to Sapling are Merging transaction and transaction assembly. The SNARKs proving time takes around two seconds for each input and one second for each output, dominating the transaction generation. The verification is noticeably slower with support for types but approximately equal to Zcash. Overall, we notice that the impact of the additional constraints required for our protocol are minimal while providing additional functionality.

## Acknowledgements

We thank the anonymous reviewers and our shepherd, Sherman Chow, for their helpful comments in improving this work. This work was partially funded by the Carlsberg Foundation under the Semper Ardens Research Project CF18-112 (BCM), the Sapere Aude: DFF-Starting Grant number 0165-00079B "Foundations of Privacy Preserving and Accountable Decentralized Protocols" and by Input Output (iohk.io) through their funding of the Edinburgh Blockchain Technology Lab.

## References

- [1] Kurt M. Alonso and Jordi Herrera Joacomartí. Monero - privacy in the blockchain. Cryptology ePrint Archive, Report 2018/535, 2018. <https://eprint.iacr.org/2018/535>.
- [2] Karim Baghery, Markulf Kohlweiss, Janno Siim, and Mikhail Volkhov. Another look at extraction and randomization of groth's zk-SNARK. Cryptology ePrint Archive, Report 2020/811, 2020. <https://eprint.iacr.org/2020/811>.
- [3] Carsten Baum, Bernardo David, and Tore Kasper Frederiksen. P2dex: Privacy-preserving decentralized cryptocurrency exchange. In Kazuo Sako and Nils Ole Tippenhauer, editors, *Applied Cryptography and Network Security*, pages 163–194, Cham, 2021. Springer International Publishing.
- [4] Eli Ben-Sasson, Alessandro Chiesa, Christina Garman, Matthew Green, Ian Miers, Eran Tromer, and Madars Virza. Zerocash: Decentralized anonymous payments from bitcoin. In *2014 IEEE Symposium on Security and Privacy*, pages 459–474. IEEE Computer Society Press, May 2014. 10.1109/SP.2014.36.
- [5] Sean Bowe, Alessandro Chiesa, Matthew Green, Ian Miers, Pratyush Mishra, and Howard Wu. ZEXE: Enabling decentralized private computation. In *2020 IEEE Symposium*



- on *Security and Privacy*, pages 947–964. IEEE Computer Society Press, May 2020. 10.1109/SP40000.2020.00050.
- [6] Benedikt Bünz, Shashank Agrawal, Mahdi Zamani, and Dan Boneh. Zether: Towards privacy in a smart contract world. In Joseph Bonneau and Nadia Heninger, editors, *FC 2020*, volume 12059 of *LNCS*, pages 423–443. Springer, Heidelberg, February 2020. 10.1007/978-3-030-51280-4\_23.
- [7] Matteo Campanelli, Felix Engelmann, and Claudio Orlandi. Zero-knowledge for homomorphic key-value commitments with applications to privacy-preserving ledgers. *Cryptology ePrint Archive*, Report 2021/1678, 2021. <https://eprint.iacr.org/2021/1678>.
- [8] Shumo Chu, Qiudong Xia, and Zhenfei Zhang. Manta: Privacy preserving decentralized exchange. *Cryptology ePrint Archive*, Report 2020/1607, 2020. <https://ia.cr/2020/1607>.
- [9] Apoorva Deshpande and Maurice Herlihy. Privacy-preserving cross-chain atomic swaps. In Matthew Bernhard, Andrea Bracciali, L. Jean Camp, Shin'ichiro Matsuo, Alana Maurushat, Peter B. Rønne, and Massimiliano Sala, editors, *FC 2020 Workshops*, volume 12063 of *LNCS*, pages 540–549. Springer, Heidelberg, February 2020. 10.1007/978-3-030-54455-3\_38.
- [10] Donghui Ding, Kang Li, Linpeng Jia, Zhongcheng Li, Jun Li, and Yi Sun. Privacy protection for blockchains with account and multi-asset model. *China Communications*, 16(6):69–79, 2019.
- [11] Felix Engelmann, Lukas Müller, Andreas Peter, Frank Kargl, and Christoph Bösch. SwapCT: Swap confidential transactions for privacy-preserving multi-token exchanges. *PoPETs*, 2021(4):270–290, October 2021. 10.2478/popets-2021-0070.
- [12] Prastudy Fauzi, Sarah Meiklejohn, Rebekah Mercer, and Claudio Orlandi. Quisquis: A new design for anonymous cryptocurrencies. In Steven D. Galbraith and Shiho Moriai, editors, *ASIACRYPT 2019, Part I*, volume 11921 of *LNCS*, pages 649–678. Springer, Heidelberg, December 2019. 10.1007/978-3-030-34578-5\_23.
- [13] Georg Fuchsbauer, Michele Orrù, and Yannick Seurin. Aggregate cash systems: A cryptographic investigation of Mimblewimble. In Yuval Ishai and Vincent Rijmen, editors, *EUROCRYPT 2019, Part I*, volume 11476 of *LNCS*, pages 657–689. Springer, Heidelberg, May 2019. 10.1007/978-3-030-17653-2\_22.
- [14] Zhimin Gao, Lei Xu, Keshav Kasichainula, Lin Chen, Bogdan Carbutar, and Weidong Shi. Private and atomic exchange of assets over zero knowledge based payment ledger. *arXiv preprint arXiv:1909.06535*, 2019.
- [15] Lorenzo Grassi, Daniel Kales, Dmitry Khovratovich, Arnab Roy, Christian Rechberger, and Markus Schofnegger. Starkad and Poseidon: New hash functions for zero knowledge proof systems. *Cryptology ePrint Archive*, Report 2019/458, 2019. <https://eprint.iacr.org/2019/458>.
- [16] Daira Hopwood, Sean Bowe, Taylor Hornby, and Nathan Wilcox. Zcash protocol specification. *GitHub: San Francisco, CA, USA*, 2016.
- [17] Thomas Kerber, Aggelos Kiayias, and Markulf Kohlweiss. Kachina—foundations of private smart contracts. In *2021 IEEE 34th Computer Security Foundations Symposium (CSF)*, pages 1–16. IEEE, 2021.
- [18] Ahmed E. Kosba, Andrew Miller, Elaine Shi, Zikai Wen, and Charalampos Papamanthou. Hawk: The blockchain model of cryptography and privacy-preserving smart contracts. In *2016 IEEE Symposium on Security and Privacy*, pages 839–858. IEEE Computer Society Press, May 2016. 10.1109/SP.2016.55.
- [19] Russell W. F. Lai, Viktoria Ronge, Tim Ruffing, Dominique Schröder, Sri Aravinda Krishnan Thyagarajan, and Jiafan Wang. Omniring: Scaling private payments without trusted setup. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *ACM CCS 2019*, pages 31–48. ACM Press, November 2019. 10.1145/3319535.3345655.
- [20] Sarah Meiklejohn and Rebekah Mercer. Möbius: Trustless tumbling for transaction privacy. *PoPETs*, 2018(2):105–121, April 2018. 10.1515/popets-2018-0015.
- [21] Andrew Poelstra, Adam Back, Mark Friedenbach, Gregory Maxwell, and Pieter Wuille. Confidential assets. In Aviv Zohar, Ittay Eyal, Vanessa Teague, Jeremy Clark, Andrea Bracciali, Federico Pintore, and Massimiliano Sala, editors, *FC 2018 Workshops*, volume 10958 of *LNCS*, pages 43–63. Springer, Heidelberg, March 2019. 10.1007/978-3-662-58820-8\_4.
- [22] Samuel Steffen, Benjamin Bichsel, Mario Gersbach, Noa Melchior, Petar Tsankov, and Martin T. Vechev. zkay: Specifying and enforcing data privacy in smart contracts. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *ACM CCS 2019*, pages 1759–1776. ACM Press, November 2019. 10.1145/3319535.3363222.
- [23] Sam M Werner, Daniel Perez, Lewis Gudgeon, Aria Klages-Mundt, Dominik Harz, and William J Knottenbelt. Sok: Decentralized finance (defi). *arXiv preprint arXiv:2101.08778*, 2021.
- [24] Jiahua Xu, Nazariy Vavryk, Krzysztof Paruch, and Simon Cousaert. Sok: Decentralized exchanges (dex) with automated market maker (amm) protocols. *arXiv preprint arXiv:2103.12732*, 2021.
- [25] Zheng Yi, Howard Ye, Patrick Dai, Sun Tongcheng, and Vladislav Gelfer. Confidential assets on MimbleWimble. *Cryptology ePrint Archive*, Report 2019/1435, 2019. <https://eprint.iacr.org/2019/1435>.

## A Proofs of Knowledge

We quickly remind the reader of the classical security properties of a NIZK: The scheme is said to be *complete*, if  $\forall(\text{stmt}, w) \in \mathcal{R}_{\mathcal{L}}$  we have

$$\Pr \left[ \begin{array}{l} (\text{crs}, \cdot) \leftarrow \text{Setup}(\mathcal{R}_{\mathcal{L}}); \\ \text{Verify}(\text{crs}, \text{stmt}, \text{Prove}(\text{crs}, \text{stmt}, w)) = 1 \end{array} \right] = 1$$

The scheme is *knowledge sound* (KS) if for all PPT  $\mathcal{A}$  there exists an extractor  $\mathcal{E}$  such that

$$\Pr \left[ \begin{array}{l} (\text{crs}, \cdot) \leftarrow \text{Setup}(\mathcal{R}_{\mathcal{L}}); \\ (\text{stmt}, \pi) \leftarrow \mathcal{A}(\text{crs}); w \leftarrow \mathcal{E}(\text{stmt}, \pi) : \\ \text{Verify}(\text{crs}, \text{stmt}, \pi) = 1 \wedge \\ (\text{stmt}, w) \notin \mathcal{R}_{\mathcal{L}} \end{array} \right] \leq \text{negl}(\lambda)$$

Finally, we say that the scheme is (computational) *zero-knowledge* if for any PPT adversary  $\mathcal{A}$  and  $\mathcal{R}_{\lambda}$ ,  $|\varepsilon_0 - \varepsilon_1| \leq \text{negl}(\lambda)$ , where

$$\varepsilon_b = \Pr \left[ (\text{crs}, \tau) \leftarrow \text{Setup}(\mathcal{R}_{\mathcal{L}}) : \mathcal{A}^{\mathcal{S}_{b, \text{crs}, \tau}}(\text{crs}) = 1 \right]$$

The oracle  $\mathcal{S}_{b, \text{crs}, \tau}$  on input  $(\text{stmt}, w)$  asserts that  $(\text{stmt}, w) \in \mathcal{R}_{\lambda}$  and then returns  $\pi = \text{Prove}(\text{crs}, \text{stmt}, w)$  if  $b = 0$ , and  $\pi = \text{Sim}(\text{crs}, \tau, \text{stmt})$  if  $b = 1$ .

We additionally require a stronger variant of the KS property called simulation-extractability, where extraction is possible in the presence of simulated proofs.

**Definition A.1** (Simulation-Extractability [2]). *We say that a non-interactive argument is weakly simulation-extractable (SE) if for any PPT adversary  $\mathcal{A}$  there exists a PT extractor  $\mathcal{E}_{\mathcal{A}}$  such that for  $\mathcal{R}_{\lambda}$ ,*

$$\Pr \left[ \begin{array}{l} (\text{crs}, \tau) \leftarrow \text{Setup}(\mathcal{R}_{\mathcal{L}}); \\ (\text{stmt}, \pi) \leftarrow \mathcal{A}^{\mathcal{S}_{\text{crs}, \tau}}(\text{crs}); \\ w \leftarrow \mathcal{E}_{\mathcal{A}}(\text{stmt}, \pi) : \\ \text{Verify}(\text{crs}, \text{stmt}, \pi) = 1 \wedge \\ (\text{stmt}, w) \notin \mathcal{R}_{\lambda} \wedge \text{stmt} \notin Q \end{array} \right] \leq \text{negl}(\lambda)$$

where  $\mathcal{S}_{\text{crs}, \tau}(\text{stmt})$  is a simulator oracle that calls  $\text{Sim}(\text{crs}, \tau, \text{stmt})$  internally, and also records  $\text{stmt}$  into  $Q$ .

The SE property guarantees that a NIZK in a certain sense “instance-binding”. E.g. if a certain game produces proofs  $\{\pi_i\}$  for a language with a hard sublanguage (e.g.  $\text{stmt} = (\text{stmt}_1, \text{stmt}_2)$  and  $\text{stmt}_1 = f(w_1)$  where  $f$  is a one-way function), we can simulate these proofs, and by SE  $\mathcal{A}$  can only return either these simulated proofs for the *same* instance, or we can extract  $w$  from them (in which case the reduction essentially inverts  $f$ ). This means that  $\mathcal{A}$  cannot change the  $\text{stmt}_2$  part in  $\{\pi_i\}$ . In practice this means that if the honest  $\pi$  hides certain secrets (a key, or a secret random value), then  $\mathcal{A}$  cannot maul  $\pi$  into proofs for a “slightly different” instance.

## B Zerocash-Style OTA

To use Zerocash style one-time-accounts, we provide an instantiation which is compatible with Zerocash.

Given a key anonymous public key encryption scheme PKE, a labeled PRF and a commitment scheme (Commit). The randomness space is  $\mathbb{S} := (\{0, 1\}^{\lambda})^3$ , the key space is  $\mathbb{K} := \{0, 1\}^{\lambda}$  and the message space consists of two  $\lambda$ -bit integers  $\mathbb{M} := \mathbb{Z}_{2^{\lambda}}^2$  the first representing an amount and the second identifying a type. The randomness space  $\Xi$  is specified by PKE.

$\mathbf{p} \leftarrow \text{Setup}(1^{\lambda})$ : Initializes the SNARK parameters  $\mathbf{p}$ .  
 $(\text{sk}, \text{pk}) \leftarrow \text{KeyGen}()$ : Samples  $a_{\text{sk}}$  randomly from  $\mathbb{K}$  and defines  $a_{\text{pk}} \leftarrow P(\text{sk}) := \text{PRF}_{a_{\text{sk}}}^{\text{addr}}(0)$ ; samples an encryption key pair  $\text{sk}_{\text{enc}}, \text{pk}_{\text{enc}}$ . Return  $\text{sk} = (a_{\text{sk}}, \text{sk}_{\text{enc}}), \text{pk} = (a_{\text{pk}}, \text{pk}_{\text{enc}})$ . The relationship between secret and public keys of the encryption keys is defined in the encryption scheme.

$\text{note} \leftarrow \text{Gen}(\text{pk}, \vec{a}, r)$ : Parse  $(\text{rk}, \text{rc}, \text{rn}) \leftarrow r$  and  $(a_{\text{pk}}, \text{pk}_{\text{enc}}) \leftarrow \text{pk}$ . Commit to  $(a_{\text{pk}}, \text{rn})$  with randomness  $\text{rk}$  as commitment  $\text{com}$  and then commit to  $\text{com}, \vec{a}$  with randomness  $\text{rc}$ . So  $\text{note} = \text{Commit}(\text{Commit}(a_{\text{pk}}, \text{rn}; \text{rk}), \vec{a}; \text{rc})$ .

$C \leftarrow \text{Enc}(\text{pk}, (\vec{a}, r), \xi)$ : Parse  $(a_{\text{pk}}, \text{pk}_{\text{enc}}) \leftarrow \text{pk}$ . Encrypt  $(\vec{a}, r)$  with  $\text{pk}_{\text{enc}}$  to ciphertext  $C$  and return  $C$ .

$(\vec{a}', r')/\perp \leftarrow \text{Receive}(\text{note}, C, \text{sk})$ : Parse  $(a_{\text{sk}}, \text{sk}_{\text{enc}}) \leftarrow \text{sk}$ . Decrypt  $C$  with  $\text{sk}_{\text{enc}}$  to  $(\vec{a}', r')$ . Parse  $(\text{rk}', \text{rc}', \text{rn}') \leftarrow r'$  and verify that these values recreate the commitment  $\text{note} = \text{Commit}(\text{Commit}(P(a_{\text{sk}}), \text{rn}'; \text{rk}'), \vec{a}'; \text{rc}')$ . If decryption fails or the commitment does not match the note, return fail.

$\text{nul} \leftarrow \text{NulEval}(\text{sk}, r)$ : Parse  $(a_{\text{sk}}, \text{sk}_{\text{enc}}) := \text{sk}$  and  $(\text{rk}, \text{rc}, \text{rn}) := r$ . Evaluate  $\text{nul} = \text{PRF}_{a_{\text{sk}}}^{\text{sn}}(\text{rn})$  and return  $\text{nul}$ .

**Theorem B.1** (OTA Binding). *When using a binding commitment scheme Commit, the Zerocash commitment scheme is binding according to Definition 3.3*

*Proof.* Let  $\mathcal{A}$  return  $\text{pk}_0, \vec{a}_0, r_0, \text{pk}_1, \vec{a}_1, r_1$  such that  $\text{Gen}(\text{pk}_0, \vec{a}_0, r_0) = \text{Gen}(\text{pk}_1, \vec{a}_1, r_1)$  and  $\vec{a}_0 \neq \vec{a}_1$ . Parsing  $\forall i \in \{0, 1\} : (\text{rk}_i, \text{rc}_i, \text{rn}_i) \leftarrow r_i$ , the first condition implies the equality  $\text{Commit}(\text{Commit}(\text{pk}_0, \text{rn}_0; \text{rk}_0), \vec{a}_0; \text{rc}_0) = \text{Commit}(\text{Commit}(\text{pk}_1, \text{rn}_1; \text{rk}_1), \vec{a}_1; \text{rc}_1)$  but the commitments have different values for  $\vec{a}$ . This breaks the binding property of the commitment scheme.  $\square$

**Theorem B.2** (Note Uniqueness). *Zerocash OTA style notes are unique according to Definition 3.5*

*Proof.* Trivially follows from binding of the underlying commitment scheme – since both  $\text{pk}$  and  $\vec{a}$  are commitment messages, finding a collision (even with adversar-

ially chosen randomness) amounts to breaking commitment binding.  $\square$

**Theorem B.3** (Nullifier Uniqueness). *The Zerocash OTA scheme has unique nullifiers (Definition 3.7), if the PRF is secure and the commitment scheme is binding.*

*Proof.* Let  $\mathcal{A}$  return  $(\text{sk}_0, r_0, \vec{a}_0, \text{sk}_1, r_1, \vec{a}_1)$  which generates the same note but different nullifiers  $\text{nul}_0$  and  $\text{nul}_1$ . The binding commitment scheme ensures that  $(\text{sk}_0, r_0, \vec{a}_0) = (\text{sk}_1, r_1, \vec{a}_1)$ . As  $\text{NulEval}$  is deterministic, the nullifiers are equal, contradicting our assumption.  $\square$

**Theorem B.4** (Nullifier Pseudorandomness). *The Zerocash OTA nullifiers are pseudorandom (Definition 3.6), if the PRF is secure and the commitment scheme is hiding.*

*Proof.* The pseudorandomness experiment is exactly the same as standard pseudorandomness, except  $\mathcal{A}$  is given (1) a public key, (2) an oracle access to  $\text{Receive}(\cdot, \cdot, \text{sk})$ . The public key is computed as  $\text{PRF}_{a_{\text{sk}}}^{\text{addr}}(0)$ , and so by pseudorandomness of this PRF (with a different domain) can be replaced by a random value  $\psi$ , so knowledge of an extra random value does not help  $\mathcal{A}$  to distinguish. Regarding the Receive oracle, it performs decryption with an unrelated  $\text{sk}_{\text{enc}}$ , so it does not interfere with the main pseudorandomness reduction, since this oracle does not use  $a_{\text{sk}}$  (it uses  $\psi$  which we already argued to be random and thus irrelevant).  $\square$

**Theorem B.5** (OTA Privacy). *Using a hiding commitment scheme and a IND-CCA, key anonymous (IK-IND) encryption scheme, the Zerocash-style OTA is private according to Definition 3.4*

*Proof.* This property trivially reduces to the assumptions mentioned, since its game is essentially a merger of IND-CCA, IK-IND, and a commitment hiding, where all the data is showed simultaneously, depending on  $b$ . Formally, we

1. switch nullifiers to random values,
2. switch the commitments from  $b = 0$  to  $b = 1$ , and then
3. switch messages by IND-CCA, and (4) keys by IK-IND.

In the steps (1) and (2) we can still simulate  $\mathcal{O}_{\text{Rcv}}$  since in both cases it uses a keypair that is unrelated to the

PRF keypair. In (3) and (4) we just wire the challenge oracles through to  $\mathcal{O}_{\text{Rcv}}$ , since  $\mathcal{O}_{\text{Rcv}}$  is essentially a CCA-like oracle.  $\square$

## C Proof of HID-OR for Sparse Pedersen Commitment

*Proof of Lemma 2.1.* Let the hash function be modelled by the random oracle, and  $H(\text{ty}_i) = G^{t_i}$ , and thus the challenger knows all  $t_i$ . The form of commitments that are given to the adversary is thus:

$$[t_b a_b + r], [t'_b a'_b + r'], \text{rc} = r - r'$$

The initial game is HID-OR<sup>0</sup>, in which  $\mathcal{A}$  sees the previous equation for  $b = 0$ . In Game<sub>1</sub> we sample  $\hat{r}$  and instead give  $\mathcal{A}$  the following:

$$[\hat{r}], [t'_0 a'_0 + r'], \text{rc} = (\hat{r} - t_0 a_0) - r'$$

Note that  $t_0 a_0 + r$  and  $\hat{r}$  are both uniform, so  $\mathcal{A}$  does not see the difference, the transition is perfect. Similarly we equivocate the second commitment:

$$[\hat{r}], [\hat{r}'], \text{rc} = (\hat{r} - \hat{r}') - (t_0 a_0 - t'_0 a'_0)$$

And now we “switch” the elements to the case of  $b = 1$ . Let  $T = \{t_0, t'_0\} = \{t_1, t'_1\}$ . Given  $\Delta_{\text{ty},0} = \Delta_{\text{ty},1}$  for all  $\text{ty} \in T$  we have:

$$t_0 a_0 - t'_0 a'_0 = \sum_{t \in T} t \Delta_{\text{ty},0} = \sum_{t \in T} t \Delta_{\text{ty},1} = t_1 a_1 - t'_1 a'_1$$

Therefore what  $\mathcal{A}$  sees now is:

$$[\hat{r}], [\hat{r}'], \text{rc} = (\hat{r} - \hat{r}') - (t_1 a_1 - t'_1 a'_1)$$

So we can proceed with replacing  $\hat{r}$  and  $\hat{r}'$  “back” into real commitments to the  $b = 1$  values similarly to how we abstracted them in the first steps of the proof. The end result is a distribution that  $\mathcal{A}$  sees for  $b = 1$  (HID-OR<sup>1</sup>):  $[t_1 a_1 + r], [t'_1 a'_1 + r'], \text{rc} = r - r'$  so our hiding holds with probability 1 (is perfect).  $\square$

**Corollary C.1.** *Assuming commitment hiding, Lemma 2.1 holds even if  $\mathcal{A}$  provides two sets of size  $n$  (not just pairs) of input and output commitments,  $\{(\text{ty}_{i,t}, a_{i,t}), (\text{ty}'_{i,t}, a'_{i,t})\}_{i,t=0,0}^{n,1}$ , as long as they still jointly sum to the same values per type.*

*Proof.* The proof is exactly the same as the previous one, except that we first “idealise”  $2n$  commitments (and not just two) from  $b = 0$ , and then de-idealise them all back. The transition logic in the middle holds similarly because  $\Delta_{\text{ty},0} = \Delta_{\text{ty},1}$ .  $\square$

## D Privacy Proof

*Proof of Theorem 6.2, deferred reduction.* Formally, we construct a series of games transforming  $\text{Privacy}_{\mathcal{A}}^0(1^\lambda)$  into  $\text{Privacy}_{\mathcal{A}}^1(1^\lambda)$ .

**Game<sub>1</sub>:** Replace all honest NIZKs (created inside  $\mathcal{O}_{\text{Spend}}$  and  $\text{SignTx}$ ) to simulated NIZKs.

**Game<sub>2</sub>:** By nullifier pseudorandomness, replace all the honest nullifiers  $\mathcal{A}$  has not seen through  $\mathcal{O}_{\text{Spend}}$  by the nullifiers over a different set of secret keys.

To switch the  $\text{NulEval}(\text{sk}_0, r_0)$  to  $\text{NulEval}(\text{sk}_1, r_1)$ , we first pick a random function  $f$  and switch *all* evaluations of  $\text{NulEval}(\text{sk}_0, r)$  to  $f(r)$ . Since the nullifiers in  $\text{tx}^{\mathcal{H}}$  have not been queried previously in  $\mathcal{O}_{\text{Spend}}$ , the evaluation of  $f(r_0)$  is the first one in the game. Hence, it is equivalent to returning a random value  $\psi$  instead of  $f(r_0)$ . Then, again by pseudorandomness, we return all the other (not related to  $\text{tx}^{\mathcal{H}}$ ) evaluations of  $f(\cdot)$  back to  $\text{NulEval}(\text{sk}_0, \cdot)$ . We now perform exactly the same steps then for the second key, de-idealizing  $\psi$  into  $\text{NulEval}(\text{sk}_1, r_1)$ . First we replace all evaluations of  $\text{NulEval}(\text{sk}_1, r)$  to  $f(r)$ , then we observe that  $r_1$  has not yet been queried to  $f(\cdot)$ , so it is equivalent to return  $f(r_1)$  instead of  $\psi$ . Then we replace all  $f(r)$  back to  $\text{NulEval}(\text{sk}_1, r)$ , *including* the target nullifier in  $\text{tx}^{\mathcal{H}}$ . This is how we have just replaced  $\text{NulEval}(\text{sk}_0, r_0)$  by  $\text{NulEval}(\text{sk}_1, r_1)$ . Repeat the procedure for all pairs of nullifiers in  $\text{tx}^{\mathcal{H}}$  in any order (i.e. it does not matter what the source and target of the replacement is since all values are pseudorandom).

**Game<sub>3</sub>:** Replace output notes and related ciphertexts to honest secret keys by the notes and ciphertexts for  $T_1$ . Recall that output notes to  $\mathcal{A}$  keys must have exactly the same inputs in both cases, so their distribution is exactly equal. This step relies on the OTA privacy and again on nullifier pseudorandomness.

By OTA privacy one can replace  $\text{OTA.Gen}(\text{pk}_b, (\text{ty}_b, v_b), r)$  and  $\text{OTA.Enc}(\text{pk}_b, (\text{ty}_b, v_b, r), \xi)$  from  $b = 0$  to  $b = 1$  even if there are Receive calls in the game before and after the replacement. (After, we cannot query the challenge note in a CCA fashion, which is irrelevant since we replace it in the very end of the game.)

The only detail left now is that there are also nullifier evaluations that use the challenge  $\text{sk}_b$ , but by nullifier pseudorandomness we can replace these nullifiers by consistent random values before the OTA privacy switch (we can do that in the presence of Receive evaluations because of the Receive oracle in the pseudorandomness definition), perform the privacy switch, and then return the real  $\text{NulEval}$  evaluations back.

**Game<sub>4</sub>:** Replace intermediary commitments com with the values from  $T_1$ . Here we apply HID-OR for vectors, as described in Corollary C.1: replacing one set of typed commitments with another, given that  $\Delta_{\text{ty},0} = \Delta_{\text{ty},1}$  for each type. This transition is perfect.

**Game<sub>5</sub>:** Remove the simulation, create real proofs according to the new data in  $T_1$ . From the first game, we have not modified anything outside of the scope of  $\text{EvalTree}$ , so all NIZKs in  $\mathcal{O}_{\text{Spend}}$  queries are “restored” from simulations without any issue, since they are for exactly the same data. Now it is a matter of a completeness check to make sure that the new transaction constructed in  $\text{EvalTree}$  has the same distribution as  $\text{tx}_{\mathcal{H},1}$  except for the (yet simulated) proofs, and in particular it satisfies  $\mathcal{L}^{\text{Spend}}$  and  $\mathcal{L}^{\text{Output}}$ . So now we enable real NIZKs back as well, and by zero-knowledge we obtain honest proofs for  $T_1$ . Now  $\text{Game}_5$  is equivalent to  $\text{Privacy}_{\mathcal{A}}^1(1^\lambda)$ , which concludes the privacy proof.

## E Balance Proof

*Proof of Theorem 6.3.* We start with the inductive step first, proving the single-tx **Balance-Single** property, presented next, which asserts that each separate transaction is well-formed (with respect to the data extracted from NIZKs), and in particular, balanced. Additionally, it allows introducing the extractor into the balance game, which we will use later.

### Balance-Single $_{\mathcal{A}, \mathcal{E}_{\mathcal{A}}}(1^\lambda)$

```

p ← Setup(1λ)
(st, tx, σ) ← A(p)
Parse tx as ( {nuliT }i=1|S| { (noteiT, Ci) }i=1|T|, {Δa:ty}ty∈Ty )
ptxε@(Sε, Tε) ← EA(p, tx, σ)
Parse Sε as { (skiS, noteiS, nuli, pathi, (aiS, tyiS, riS) }i=1|S|
Parse Tε as { (pkiT, noteiT, CiT, (aiT, tyiT, riT) }i=1|T|
return Verify(st, tx, σ) = 1 ∧
      (CheckPTx(ptxε, st) ≠ 1 ∨ CompleteTx(ptxε) ≠ tx)

```

Formally, we claim that for every  $\mathcal{A}$  exists  $\mathcal{E}_{\mathcal{A}}$  s.t.  $\Pr[\text{Balance-Single}_{\mathcal{A}, \mathcal{E}_{\mathcal{A}}}(1^\lambda) = 1] \leq \text{negl}(\lambda)$ .

*Proof (Balance-Single).* We follow the  $\text{CheckPTx}$  code and show why it succeeds. The check that  $\text{nul}_i \notin \text{st.MT}$  is directly checked in  $\text{Verify}$ . All except the last (balance) checks reduce to the NIZK knowledge-soundness.

Looking at every  $\mathcal{L}^{\text{Spend}}$  NIZK first, the fact it verifies implies that we extract  $w = (\text{path}_i, \text{sk}_i^S, a_i^S, \text{ty}_i^S, r_i^S, \text{rc}_i^S)$ , which by NIZK soundness guarantees for all  $i \in [|\mathcal{S}|]$ :

$$\begin{aligned} \text{st.MT}[\text{path}_i] &= \text{note}_i^S \wedge \\ \text{nul}_i &= \text{OTA.NulEval}(\text{sk}_i^S, r_i^S) \wedge \\ \text{note}_i^S &= \text{OTA.Gen}(\text{OTA.P}(\text{sk}_i), (a_i^S, \text{ty}_i^S), r_i^S) \end{aligned}$$

The (1) and (3) are by  $\mathcal{L}^{\text{open}}$ , and (2) is by  $\mathcal{L}^{\text{nul}}$  (see  $\mathcal{L}^{\text{spend}}$  structure); Similarly, the output NIZKs guarantee (through  $\mathcal{L}^{\text{open}}$ ) the existence of  $(\text{pk}^T, a^T, \text{ty}^T, r^T, \text{rc}^T)$  s.t.:

$$\forall i \in [|\mathcal{T}|] : \text{note}_i^T = \text{OTA.Gen}(\text{pk}_i^T, (a_i^T, \text{ty}_i^T), r_i^T)$$

Which is again by  $\mathcal{L}^{\text{open}}$ . The last (balance) check that  $\forall \text{ty} \in \{\text{ty}_i^S\}_{i=1}^{|\mathcal{S}|} \cup \{\text{ty}_i^T\}_{i=1}^{|\mathcal{T}|} : \sum_{a \in \{a_i^S | \text{ty}_i^S = \text{ty}\}_{i=1}^{|\mathcal{S}|}} a - \sum_{a \in \{a_i^T | \text{ty}_i^T = \text{ty}\}_{i=1}^{|\mathcal{T}|}} a = \Delta_{\text{ty}}$  succeeds by the commitment homomorphic property in  $\text{Verify}(\cdot, \cdot, \cdot)$ . We know by the previous steps that  $a_i$  are the same values that are committed into  $\text{com}_i$ , and the commitments sum to the identity, which means that the committed values  $a_i$  sum to 0.  $\square$

The main balance property follows inductively by following the  $(\text{st}, \text{st}', \text{tx}, \sigma) \in \text{Ins}'$  loop step-by-step. But before we present the main proof reasoning, we must introduce several auxiliary constructions into the game. At each step we will trace the notes that  $\mathcal{A}$  can spend, and will argue that both during the game, and in the end of it, the sum of these is not enough to break the balance predicate.

- Instead of computing  $v_{\mathcal{A}}$  in a separate loop before the main loop, we will move it *into* the main loop. Now,  $v_{\mathcal{A}}$  is updated whenever the balance game locates a note  $\in I$  in the outputs of  $\text{tx}$  (this note is still attempted to be received by  $\text{SK}^*$ ). This allows us to track the successive stake accumulation of  $\mathcal{A}$ , but only on those coins that it claims in  $I$ .
- We add a set<sup>11</sup>  $N_{\overline{\mathcal{H}}}$ , initially populated with notes in  $I_0$ . Each loop iteration will
  1. remove adversarial input notes from it, and
  2. add to it all the output notes that cannot be received by  $\text{SK}$ .

For both actions we need the data extracted by  $\text{Balance-Single}$ , that “reveals” the notes behind the nullifiers, including adversarial ones. The new variable  $v_{\overline{\mathcal{H}}}$  will track balances inside  $N_{\overline{\mathcal{H}}}$ , and is updated whenever  $N_{\overline{\mathcal{H}}}$  is changed.  $N_{\overline{\mathcal{H}}}$  and  $v_{\overline{\mathcal{H}}}$  track *all* the transient adversarial transactions, not counted

in the final  $v_{\mathcal{A}}$  (which only sums up the “resulting” assets of  $\mathcal{A}$ ). E.g.  $\mathcal{A}$  can “refresh” the note  $\text{note} \in I_0$  by moving its funds to another adversarial note'  $\notin I$ , and without  $N_{\overline{\mathcal{H}}}$  this will not be reflected anywhere.

- We insert the assertion  $\exists \text{ty} \in T : v_{\overline{\mathcal{H}}}[\text{ty}] \leq v_0[\text{ty}] + v_{\mathcal{H}-}[\text{ty}] - v_{\mathcal{H}+}[\text{ty}]$  into the end of each loop. Clearly this only increases adversarial chances to win by providing an extra winning condition.

For example, if this predicate fails in the middle of the loop w.r.t. some intermediate  $\text{st}'$ ,  $\mathcal{A}$  could as well choose  $\text{st}'$  for its challenge response, and win the original game with the assertion in the end. And also clearly using  $v_{\overline{\mathcal{H}}}$  instead of  $v_{\mathcal{A}}$  in the predicate only gives  $\mathcal{A}$  more advantage, since in the end  $v_{\mathcal{A}} \leq v_{\overline{\mathcal{H}}}$ : by anti-theft,  $\mathcal{A}$  cannot spend honest nullifiers, and thus it can only claim notes from  $N_{\overline{\mathcal{H}}}$  (and can only do it once by nullifier uniqueness and OTA binding), but  $v_{\mathcal{A}}$  is calculated w.r.t.  $I \subset N_{\overline{\mathcal{H}}}$ .

In the base case of the induction, before the first loop (before any transaction is processed), the predicate holds, since  $v_{\overline{\mathcal{H}}}[\text{ty}] = v_0[\text{ty}]$  and  $v_{\mathcal{H}+} = v_{\mathcal{H}-} = 0$ .

Assume now that the predicate holds in the beginning of the loop, we will show it persists through the loop execution. Fix a type  $\text{ty}$ , and apply the following reasoning to each type in the transaction. Compute the local values  $v_{\overline{\mathcal{H}-}}^{\Delta}, v_{\overline{\mathcal{H}+}}^{\Delta}$  as prescribed by the game. Compute  $v_{\mathcal{A}}^{\Delta}$  by looking at all the outputs in  $\text{tx}_{\mathcal{A}}$ , and summing the corresponding values that can be received using  $\text{SK}^*$ . These values are the same as the values extracted by a single-step definition, by OTA output note binding.

Compute  $v_{\overline{\mathcal{H}}}^{\Delta}$  from the extracted (input and output) NIZKs by adding all the adversarial output values and subtracting all the adversarial input values.

Now we need to prove that  $v_{\overline{\mathcal{H}}}^{\Delta} \leq v_{\overline{\mathcal{H}-}}^{\Delta} - v_{\overline{\mathcal{H}+}}^{\Delta}$  ( $v_0$  does not ever change). But this reduces to the commitment balancing condition:  $v_{\overline{\mathcal{H}}}^{\Delta} = \sum a_i^{\mathcal{T}\mathcal{A}} - \sum a_i^{\mathcal{S}\mathcal{A}}$ , while  $v_{\mathcal{H}-} - v_{\mathcal{H}+} = \sum a_i^{\mathcal{S}\mathcal{H}} - \sum a_i^{\mathcal{T}\mathcal{H}}$ , and by the last (balancing) check in  $\text{CheckPTx}$ , guaranteed by  $\text{Balance-Single}$ , we have

$$\sum a_i^{\mathcal{S}\mathcal{A}} + \sum a_i^{\mathcal{S}\mathcal{H}} - \sum a_i^{\mathcal{T}\mathcal{A}} - \sum a_i^{\mathcal{T}\mathcal{H}} \geq 0$$

(Greater than equal, since  $\Delta_{\text{ty}}$  can be positive.) Which translates into  $-v_{\overline{\mathcal{H}}}^{\Delta} + v_{\mathcal{H}-} - v_{\mathcal{H}+} \geq 0$ , equivalent to  $v_{\overline{\mathcal{H}}}^{\Delta} \leq v_{\mathcal{H}-} - v_{\mathcal{H}+}$  as we need. So the predicate persists through the loop iteration.

Since at each point  $v_{\overline{\mathcal{H}}} \geq v_{\mathcal{A}}$ , and the predicate  $v_{\overline{\mathcal{H}}} \leq \dots$  holds in the end of each loop, the predicate in the end with  $v_{\mathcal{A}} \leq \dots$  will also hold, and thus  $\mathcal{A}$  cannot win the balance game.

<sup>11</sup> “Not honest” in the context of notes means adversarial notes and also notes that are “burned” — that cannot be accessed by anyone.