

# A Transaction-Level Model for Blockchain Privacy

François-Xavier Wicht<sup>1</sup>[0009-0005-6090-7901], Zhipeng Wang<sup>2</sup>, Duc V. Le<sup>3</sup>, and Christian Cachin<sup>1</sup>[0000-0001-8967-9213]

<sup>1</sup> Institute of Computer Science, University of Bern  
{francois-xavier.wicht,christian.cachin}@unibe.ch

<sup>2</sup> Imperial College London  
zhipeng.wang20@imperial.ac.uk

<sup>3</sup> Visa Research\*\*  
duc.le@visa.com

**Abstract.** Considerable work explores blockchain privacy notions. Yet, it usually employs entirely different models and notations, complicating potential comparisons. In this work, we use the Transaction Directed Acyclic Graph (TDAG) and extend it to capture blockchain privacy notions (PDAG). We give consistent definitions for untraceability and unlinkability. Moreover, we specify conditions on a blockchain system to achieve each aforementioned privacy notion. Thus, we can compare the two most prominent privacy-preserving blockchains – Monero and Zcash, in terms of privacy guarantees. Finally, we unify linking heuristics from the literature with our graph notation and review a good portion of research on blockchain privacy.

**Keywords:** Blockchain privacy · Untraceability · Unlinkability.

## 1 Introduction

Public blockchains are inherently transparent and most cryptocurrencies fail to meet even the basic privacy standards. For instance, each transaction conducted on the two most prominent blockchains, Bitcoin and Ethereum, fully disclose senders, recipients and exchanged amount. Furthermore, identities are quite thinly protected behind pseudonymous addresses, which can easily be compromised by various mechanisms [5, 6, 22, 40]. Once they are, restoring privacy becomes a challenging and expensive endeavor.

Many techniques, add-on solutions, or recommended behaviors exist to increase the base level of privacy. One such behavior is first for users to own multiple addresses and second to change them frequently. Yet, such a technique implies that users' wealth is scattered over several addresses. At some point, if a user wants to operate multiple of her addresses to conduct a single transaction, such

---

\*\* The main part of the work was conducted while the author was at the University of Bern.

a pattern is noticeable [2, 28, 33]. This linking is based on the assumption that multiple coins spent together have the same ownership. However, it is possible for various parties to combine their coins and send them to whomever they wish. This is the base principle of the Bitcoin add-on privacy solution *Coinjoin* [26]. A more sophisticated mechanism exists in smart-contract-enabled blockchains such as Ethereum with *Tornado Cash* [36] and related mechanisms. They also allow different users to mix their cryptocurrency together to obscure the origin of the funds. Yet, in that scenario, individuals deposit their coins and receive a note allowing them to withdraw their funds later. Thanks to zero-knowledge proofs, the link between deposits and withdrawals is not revealed, while preventing users from withdrawing funds twice. Unfortunately, the opt-in nature of these tools does not increase the default privacy of the underlying chains. Moreover, there exist multiple ways to decrease the unlinkability properties of add-on privacy tools, as shown by several studies [18, 43]. This is also possible because, in such blockchains, no inherent mechanisms exist to hide the origin of transactions. Thus, chain analysis can easily be conducted to trace coins back to their origin.

Privacy-preserving blockchains address this issue by hiding the sender of transactions and providing confidentiality for the amounts transferred. For example, Zcash builds a so-called shielded pool where transactions are anonymized via zero-knowledge proofs [4]. In contrast, Monero expands on the CryptoNote protocol to hide the source of transactions with ring signatures [39]. Zcash and Monero are thus privacy-preserving cryptocurrencies, and they have attracted a lot of attention from the public and from research. However, despite extensive studies, some questions remain open about privacy in blockchain.

- How do the different privacy mechanisms and tools compare?
- What are the relevant structures in a blockchain and the ledger that either protect privacy or break it?
- How do different shapes of privacy manifest themselves in a cryptocurrency using a blockchain?

This work provides an answer to those questions by defining two main privacy notions, *untraceability* and *unlinkability*, using a formal model of the blockchain transaction structure. For that, we build upon a pre-existing model, the “Transaction Directed Acyclic Graph” (TDAG) defined by Cachin et al. [7], and devise the Privacy-preserving transaction DAG (PDAG) to capture additional properties such as the hiding of the transaction’s source. This model allows us to contribute to the field in several aspects.

- We provide definitions for untraceability and unlinkability, and thus lift the ambiguity between those terms.
- Also, we give a structural view over blockchain privacy.
- Furthermore, we give conditions on a system for achieving each notion.
- This allows us to compare the different blockchain systems in terms of privacy guarantees.
- Finally, thanks to our notation, we can unify analysis patterns found in the literature that are usually used to link related addresses.

We first recall the building blocks of the TDAG and present the PDAG in Section 2. Section 3 tackles the definition of untraceability and unlinkability using PDAG’s elements. In Appendix A and in the full version of the paper [45], we apply our models to add-on privacy solutions, Coinjoin and Tornado Cash, and two privacy-preserving blockchains, namely Monero and Zcash.

## 2 Privacy-Preserving Transaction Directed Acyclic Graph

This section introduces the *Privacy-Preserving Transaction Directed Acyclic Graph* (PDAG). This model is based upon the Transaction Directed Acyclic Graph (TDAG) described by Cachin et al. [7]. We therefore start by presenting the TDAG, give motivations for a privacy-preserving version, and finally define the PDAG.

### 2.1 Background: Transaction Directed Acyclic Graph

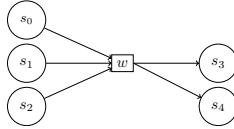
A *Transaction Directed Acyclic Graph* (TDAG) [7] is a graph  $G = (\mathcal{V}, \mathcal{E})$  with two types of vertices, namely states  $\mathcal{S}$  and witnesses  $\mathcal{W}$ , such that  $\mathcal{V} = \mathcal{S} \dot{\cup} \mathcal{W}$ . The edges, in turn, represent transitions between states and witnesses. In other words, edges  $\mathcal{E}$  represent transactional relations between states and witnesses. Edges are partitioned into three different types: consuming, observing, and producing edges, denoted respectively by,  $\mathcal{E}_C$ ,  $\mathcal{E}_O$ , and  $\mathcal{E}_P$ .

- **States** ( $\bigcirc$ ). A state  $s \in \mathcal{S}$  is the first type of vertex, denoting the output state of a transaction. It refers to an individual asset on the blockchain, a digital coin, a coin controlled by a particular address, or the state of a smart contract. A state is the result of a transaction and can transition to other states through subsequent transactions. The full context of the blockchain consists of all states that exist at a given time. A unique *genesis state*, represents the initial state of the blockchain network. There is only one genesis state, since the blockchain can only be bootstrapped once.
- **Witnesses** ( $\square$ ). A witness  $w \in \mathcal{W}$  is the second kind of vertex, representing any data in a transaction for it to be valid according to the blockchain validation rules. Every blockchain transaction contains exactly one witness.
- **Consuming edges** ( $\bigcirc \longrightarrow \square$ ). A consuming edge  $e \in \mathcal{E}_C$  connects a state  $s$  to a witness  $w$  and expresses that  $s$  is consumed by the transaction involving  $w$ . A state can be consumed at most once. After this, no other transaction may consume  $s$ . A state that is consumed has been “updated” or “overwritten” by the transaction.
- **Producing edges** ( $\square \longrightarrow \bigcirc$ ). A producing edge  $e \in \mathcal{E}_P$  connects a witness  $w$  to a state  $s$  and expresses that  $s$  is created by the transaction corresponding to  $w$ .
- **Observing edges** ( $\bigcirc \cdots \rightarrow \square$ ). An observing edge  $e \in \mathcal{E}_O$  connects a state  $s$  to a witness  $w$  and expresses that  $s$  enters into a transaction represented by  $w$ , but that it stays available to another transaction for consumption.

Those elements constitute the building blocks of a TDAG  $G$ , and note that  $G$  is bipartite. Additionally, every *transaction* in a TDAG  $G = (\mathcal{V}, \mathcal{E})$  is a subgraph denoted by  $T_w = (\mathcal{S}' \dot{\cup} \{w\}, \mathcal{E}')$  composed of a subset of vertices  $\mathcal{V}' = \mathcal{S}' \dot{\cup} \{w\} \subseteq \mathcal{V}$ , a subset of edges  $\mathcal{E}' \subseteq \mathcal{E}$ , a single witness  $w$ , a set of input states  $\mathcal{S}_I \subseteq \mathcal{S}'$  and a set of output states  $\mathcal{S}_O \subseteq \mathcal{S}'$ , such that

- Every input state in  $\mathcal{S}_I$  is a source (has indegree zero);
- Every output state in  $\mathcal{S}_O$  is a sink (has outdegree zero);
- $\mathcal{V}' = \mathcal{S}_I \dot{\cup} \mathcal{S}_O \dot{\cup} \{w\} \subseteq \mathcal{V}$ ;
- Every edge  $e$  in  $\mathcal{E}'$  is either a consuming edge  $e_C$ , or an observing edge  $e_O$  and links some input state  $s_i \in \mathcal{S}_I$  to  $w$ , or it is a producing edge and links  $w$  to some output state  $s_o \in \mathcal{S}_O$ .

In Figure 1, we depict an example of a transaction  $T_w$ . A possible interpretation of  $T_w$  is an asset transfer from three states  $s_0, s_1, s_2$  to two different states  $s_3, s_4$ . The input states of  $T_w$  are being consumed, while the output states are being produced. The produced states are now available for consumption by subsequent transactions. Regarding multiple transactions, the transaction graph must follow



**Fig. 1.** Illustration of a transaction  $T_w = (\mathcal{V}', \mathcal{E}')$ , where  $\mathcal{V}' = \{s_0, s_1, s_2, w, s_3, s_4\}$  and  $\mathcal{E}' = \{(s_0, w), (s_1, w), (s_2, w), (w, s_3), (w, s_4)\}$ . The input states  $\{s_0, s_1, s_2\}$  are consumed, whereas the output states  $\{s_3, s_4\}$  are produced.

certain rules and have a well-defined structure. The graph starts with the genesis state  $s_g$ . Regular states are produced exactly once and are consumed at most once. There are otherwise no restrictions on the number of observing edges. This construction allows the TDAG to represent UTXO-based (e.g., Bitcoin), account-based blockchains (e.g., Ethereum), and possibly others. However, we argue in the next subsection that the TDAG lacks certain elements to represent privacy-preserving blockchains (e.g., Monero and Zcash).

## 2.2 Privacy-Preserving Transactions

This section motivates and defines building blocks to introduce additional components for the TDAG to represent privacy-preserving blockchains. We motivate these additions in regard to conflict-freedom. Cachin et al. [7] give conflict-freedom conditions to a TDAG. Conflicts in a blockchain underlying a cryptocurrency occur in case of double-spending. Such a situation arises if the same unique asset is spent more than once. In UTXO-based blockchains, the system tracks unspent transaction outputs (UTXO) to determine which transactions may be

spent. On privacy-preserving chains, the UTXO set is built differently so as not to disclose any information to its users. For example, in Monero, each transaction generates a unique “key image”, and key images used more than once are rejected by the blockchain [21]. A transaction is, therefore, valid if a key image does not conflict with an existing one. Zcash’s shielded equivalent of a UTXO is a “commitment” [17]. A node can, in turn, spend a commitment by revealing a “nullifier”.

*Masking.* In a TDAG, conflict-freedom holds when every state is produced and consumed at most once, i.e., coins cannot be spent twice, nor can they be duplicated. Now, the same rule applies to privacy-preserving blockchains. Yet, state consumption (coin spending) is not directly visible to the network beyond the sender and receiver. Say, for example, that Alice mixes her input with ten other decoy inputs using a cryptographic mechanism for her to remain anonymous. The edges from these ten decoys are observing edges, since we do not alter them, but only read them. However, Alice’s spent coins are effectively consumed but appear only as being observed for non-involved parties. This introduces a new type of edges, *masking edges*  $\mathcal{E}_M$ .

- **Masking edges** ( $\circ \rightsquigarrow \square$ ). A masking edge  $e \in \mathcal{E}_M$  connects a state to a witness and is a consuming edge that hides behind cryptographic mechanisms. More precisely, it appears as an observing edge to nodes aside from the author, but is a consuming edge. As its name suggests, it masks its actual behavior. It otherwise follows the same rules as consuming edges.

Moreover, such addition brings ambiguity towards observing edges, since an outside observer may not distinguish an observing edge from a masking one. We thus call *ambiguous edges* the union of observing and masking edges, i.e.,  $\mathcal{E}_O \dot{\cup} \mathcal{E}_M$ .

*Nullifiers.* Privacy-preserving blockchains also output additional components at each transaction. These components allow constructing a cryptographic mapping with unspent transactions. This serves to prevent double-spending while hiding the source of a transaction. To be more general, hidden behaviors in the blockchain require users to output a specific value as a *nullifier*. In addition, the exact same behavior must produce the same value. Therefore, this redundancy is visible if users try to benefit from the same asset or service multiple times. The blockchain validation rules will thus forbid one such transaction since two values collide. This implies that each transaction must verify that the new values do not conflict with pre-existing ones. In turn, the total set of generated values cannot be discarded at any point, since it must be read for each transaction to prevent conflicts. We therefore introduce this additional type of state that can only be observed or produced. Depending on the implementation, this component may have different names, e.g., key images for Monero, nullifiers for Zcash, or withdrawal commitments for Tornado Cash [36]. We thus introduce a new type of vertex that we call (without loss of generality) *nullifiers*  $\mathcal{N}$ .

- **Nullifiers** ( $\diamond$ ). A nullifier  $n \in \mathcal{N}$  is a state that encapsulates cryptographic commitments which prevent users from benefiting from the same asset or service multiple times. It may only be produced and observed, but not consumed.

Those two additions lead us to call this model the nullifier-based one, as opposed to the “UTXO-based” or “account-based” model. We formalize the different components in the next subsection.

### 2.3 Privacy-Preserving Transaction DAG

We here define our model, the *Privacy-preserving Transaction DAG* (PDAG). We begin by formally specifying the PDAG and its elements, and next define a transaction in a PDAG. We use the notation  $\mathcal{E} \subseteq X \times Y$  to denote a relation  $\mathcal{E}$  between  $X$  and  $Y$ . Also, the expression  $(x, y) \in \mathcal{E}$  can be written as  $x\mathcal{E}y$ , and we define the set  $x\mathcal{E}\star$  as  $\{y : x\mathcal{E}y\}$ , and write its cardinality as  $|x\mathcal{E}\star|$ . We speak of the closed neighborhood of a vertex  $v$  to mention any adjacent vertex to  $v$ , including  $v$  itself, and denote it by  $N[v]$ .

**Definition 1 (PDAG).** *A Privacy-preserving Transactional DAG (PDAG) is a directed unweighted graph  $G = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \mathcal{S} \dot{\cup} \mathcal{N} \dot{\cup} \mathcal{W}$  are the vertices and  $\mathcal{E} = \mathcal{E}_C \dot{\cup} \mathcal{E}_M \dot{\cup} \mathcal{E}_O \dot{\cup} \mathcal{E}_P$  are the edges, where vertices and edges are each partitioned into respective sets. The sets  $\mathcal{S}, \mathcal{N}, \mathcal{W}, \mathcal{E}_C, \mathcal{E}_M, \mathcal{E}_O, \mathcal{E}_P$  constitute the components of  $G$ . The set  $\mathcal{S}$  denotes the states and contains a special state  $s_g$  called genesis. The nullifiers  $\mathcal{N}$  are states that encapsulate cryptographic commitments. The set  $\mathcal{W}$  denotes the witnesses. Edges are partitioned into four subsets, where  $\mathcal{E}_C, \mathcal{E}_M, \mathcal{E}_O \subseteq \mathcal{S} \times \mathcal{W}$  and  $\mathcal{E}_P \subseteq \mathcal{W} \times \mathcal{S}$ . It satisfies the following conditions:*

1.  $s_g$  does not have any producing or observing edges, and it has a single consuming edge:  $|\star\mathcal{E}_P s_g| = 0 \wedge |s_g\mathcal{E}_C\star| = 0 \wedge \exists! w \in \mathcal{W} : s_g\mathcal{E}_C w$ .
2. Every state, except for the genesis state, has exactly one producing edge:  $\forall s \in \mathcal{S} \setminus \{s_g\} \exists! w \in \mathcal{W} : w\mathcal{E}_P s$ .
3. Every state, except for the genesis state, may have multiple successors, but at most one among them is connected with a consuming or masking edge:  $\forall s \in \mathcal{S} : |s\mathcal{E}_C\star \dot{\cup} s\mathcal{E}_M\star| \leq 1$ .
4. Nullifiers  $n \in \mathcal{N}$  have no consuming edges and no masking edges:  $\forall n \in \mathcal{N} |n\mathcal{E}_C\star| = 0 \wedge |n\mathcal{E}_M\star| = 0$ .
5.  $G$  is weakly connected.
6.  $G$  has no cycles.

Note that for each witness  $w$ , there is a corresponding transaction  $t$ . We thus define a transaction accordingly. We use the notation  $x \prec y$  to say that  $x$  comes before  $y$  in the topological order of the graph, and call  $x$  a predecessor of  $y$ .

**Definition 2 (Transaction in a PDAG).** *Given a PDAG  $G = (\mathcal{S} \dot{\cup} \mathcal{N} \dot{\cup} \mathcal{W}, \mathcal{E})$  and a witness  $w \in \mathcal{W}$ , the transaction with witness  $w$  is the unique subgraph  $T_w = (\mathcal{S}' \dot{\cup} \mathcal{N}' \dot{\cup} \{w\}, \mathcal{E}')$   $\subset G$ , where*

- $\mathcal{S}'$  is the set of states connected to  $w$ :  $\mathcal{S} = \{s \in \mathcal{S} : s\mathcal{E}_C w \vee s\mathcal{E}_M w \vee s\mathcal{E}_O w \vee w\mathcal{E}_P s\}$ ;
- $\mathcal{N}'$  is the subset of all nullifiers produced by predecessors and the ones produced by  $T_w$ :  $\mathcal{N}' = \{w' \in \mathcal{W}, n \in \mathcal{N} : (w'\mathcal{E}_P n \wedge w' \prec w) \wedge w\mathcal{E}_P n\}$ ;
- $w \in \mathcal{W}$  is the witness of the transaction; and
- $\mathcal{E}'$  are the edges with both endpoints in  $\mathcal{S}' \dot{\cup} \{w\}$ .

Furthermore, a transaction may be the one of multiple types, depending on the number of its inputs and outputs:

- **SISO**. Single-input, single-output transactions consist of one consuming (or masking) edge that connects an input state to a witness  $w$  and of a producing edge that links  $w$  to an output state.
- **SIMO**. Single-input, and multiple-output transactions possess producing edges from  $w$  to multiple output states.
- **MISO**. Multi-input, single-output transactions have a set of multiple consuming (or masking) and observing edges.
- **MIMO**. Multi-input, multi-output transactions contain multiple consuming (or masking) and observing edges, as well as multiple producing edges.

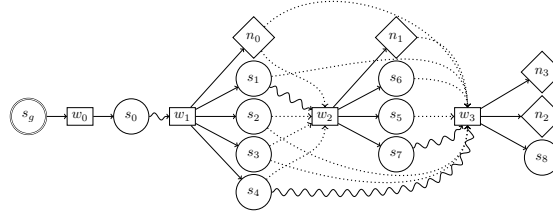
One special transaction, the initialization transaction, consists of a consuming edge linking the genesis state and a set of producing edges. The next section provides a minimal example that presents different types of transactions.

## 2.4 Example

We here assume a hypothesized privacy-preserving blockchain that uses an abstract cryptographic mechanism to hide the source of each transaction. This mechanism takes any previously produced state as decoy to the real input. For simplicity, we assume this blockchain does not incur fees for its transactions, and that the first transaction consists of a coin minting process. Each transaction outputs nullifiers to prevent double-spending. Assets are held by public keys (addresses) and to transfer said asset, one needs to prove knowledge of the private key with a signature.

In Figure 2, we depict the PDAG of a sample execution of such blockchain. The first transaction  $T_{w_0}$  consists of the minting of coins to the address of the miner ( $s_0$ ). The second transaction  $T_{w_1}$  represents the transfer from  $s_0$  to four different addresses:  $s_1, s_2, s_3, s_4$ . This transaction links  $s_0$  to  $w_1$  with a masking edge, but does not include any decoy inputs. The witness  $w_1$  verifies the validity of the signature and of the produced nullifier  $n_0$ . The third transaction  $T_{w_2}$  uses three decoy inputs  $s_2, s_3, s_4$  to transfer assets from  $s_1$  to three different addresses:  $s_5, s_6, s_7$ . These decoy keys are linked to  $w_2$  with observing edges, since they are only read (and not consumed) by the transaction. The witness  $w_2$  verifies the signature and the conflict-freedom between the previously generated nullifier ( $n_0$ ) and the newly produced one ( $n_1$ ). The last transaction  $T_{w_3}$  merges coins from two different addresses ( $s_4, s_7$ ) to a single one ( $s_8$ ). Transaction  $T_{w_3}$  takes five decoy addresses as inputs,  $s_1, s_2, s_3, s_5, s_6$ , which are linked with observing

edges. The verification of the signatures and the conflict-freedom between the previous nullifiers ( $n_0, n_1$ ) and the new one ( $n_2, n_3$ ) take place at the witness  $w_3$ .



**Fig. 2.** Illustration of a PDAG with four transactions:  $T_{w_0}$ ,  $T_{w_1}$ ,  $T_{w_2}$ , and  $T_{w_3}$ . Each transaction  $T_i$  is the closed neighborhood subgraph  $N_G[w_i] \subset G$ .

### 3 Privacy Notions

This section formalizes two privacy notions in blockchain systems with respect to the PDAG, namely untraceability and unlinkability. We first describe the adversarial model, and secondly formalize the different privacy notions.

#### 3.1 Adversarial Model

We use  $\Pi$  to denote the analyzed blockchain system and  $\mathcal{A}$  for an adversary. We assume idealized cryptographic primitives and a computationally-unbounded  $\mathcal{A}$ . However, we assume that the adversary’s focus lies on the transaction layer, and that  $\mathcal{A}$  operates passively, solely observing the transactions to compromise the privacy guarantees of  $\Pi$ . To establish the concepts of untraceability and unlinkability, we adopt the notion of anonymity sets, as introduced by Chaum [8]. Specifically, we refer to *untraceability sets* and *unlinkability sets*, respectively, which we define and quantify based on the PDAG.

#### 3.2 Untraceability

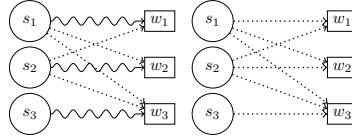
Untraceability is commonly understood in the blockchain transaction layer as the property that for each transaction, all possible senders are equiprobable [39, 24, 32, 46, 17]. In other fields of computer science, this notion is sometimes addressed as *sender anonymity* [37]. According to the PDAG, a sender is an input state linked to a witness with a consuming or with a masking edge. In turn, a sender gets a degree of untraceability if the witness of a transaction has no consuming edge, and if the number of ambiguous edges is greater than one. Otherwise, the sender of the transaction is visible or, readily guessable. We call



such a transaction trivially traceable. The input states of a transaction linked with ambiguous edges, thus form the untraceability set of the transaction, and any member of this set is equiprobable to be consumed. For  $\mathcal{A}$ , it is a matter of lifting the ambiguity behind these edges and to guess which of the states is in fact consumed. We adopt a combinatorial view of this problem and model it as a bipartite graph. This approach is similar to the study of anonymous communications introduced by Edman et al. [11] and later refined by Gierlichs et al. [14]. Given a PDAG  $G = (\mathcal{S} \dot{\cup} \mathcal{W} \dot{\cup} \mathcal{N}, \mathcal{E})$ , one can transform it into another bipartite graph  $G^* = (\mathcal{S}^* \dot{\cup} \mathcal{W}^*, \mathcal{E}^*)$  by pruning vertices and edges. We only include in  $G^*$ , witnesses of regular transactions that exchange assets between addresses, and thus exclude mining, minting or any ordering transactions, and also consider input states to these witnesses with their corresponding edges, i.e.,

- $\mathcal{W}^* = \{w^* \in \mathcal{W} \mid w^* \text{ is a regular transaction} \} \subseteq \mathcal{W}$ ,
- $\mathcal{S}^* = \{s^* \in \mathcal{S} \mid (s^*, w^*) \in \mathcal{S} \times \mathcal{W}^*\} \subseteq \mathcal{S}$ , and
- $\mathcal{E}^* = \mathcal{S}^* \times \mathcal{W}^* \subseteq \mathcal{S} \times \mathcal{W}$ .

One such bipartite graph composed of three transactions is presented in Figure 3 with  $\mathcal{S}^*$  on the left and  $\mathcal{W}^*$  on the right. Each transaction  $T_{w_i}^*$  is the closed neighborhood  $N_{G^*}[w_i]$  with  $1 \leq i \leq 2$ . If we consider the single transactions



**Fig. 3.** Three transactions  $T_{w_1}^*$ ,  $T_{w_2}^*$ , and  $T_{w_3}^*$  have a total of three input states  $s_1$ ,  $s_2$  and  $s_3$ . On the left is the graph that represents the ground-truth with masking edges, and on the right, the view of the adversary. The masking edges highlight the maximum matching that  $\mathcal{A}$  must guess.

$T_{w_1}^*$ ,  $T_{w_2}^*$ , and  $T_{w_3}^*$ , the traceability sets are  $\{s_1, s_2\}$ ,  $\{s_1, s_2\}$ , and  $\{s_1, s_2, s_3\}$ , respectively. The senders of transactions are untraceable among two for  $w_1$  and  $w_2$ , and among three for  $w_3$ . Furthermore, the masking edges form a maximum matching over the two sets of vertices. In this example, the matching is also perfect since each vertex is covered by the matching. However, this is not necessarily the case, since non-consumed states may also be part of the untraceability set. In addition, this maximum matching is usually not unique. In Figure 3, we find two different maximum matchings:  $M_1 = \{(s_1, w_1), (s_2, w_2), (s_3, w_3)\}$  and  $M_2 = \{(s_1, w_2), (s_2, w_1), (s_3, w_3)\}$ . Only  $M_1$  highlights the actual senders of transactions in this case. If this matching was unique, an adversary could trivially know each source of transactions irrespective of the incoming degree of their witness. It indeed follows from the graph structure that each witness consumes at least one state, and hence that the state consumption is represented by a maximum matching. Therefore, one way of reducing the untraceability set

is to only consider states linked with edges part of a maximum matching. This method has recently been studied by Vijayakumaran [42] and Egger et al. [12] to perform graph-based deanonymization on the CryptoNote protocol. They use the union of maximum matchings known as the Dulmage-Mendelsohn core [10] to prune edges from the bipartite graph and thus to compromise the untraceability of transactions. We here adapt the definition of the core of a bipartite graph  $G^*$  to our notation.

**Definition 3 (Core [10]).** *The core of a bipartite graph  $G^* = (\mathcal{S}^* \dot{\cup} \mathcal{W}^*, \mathcal{E}^*)$ , denoted by  $\mathbf{core}(G^*) = (\mathcal{S}^* \dot{\cup} \mathcal{W}^*, \mathcal{E}')$ , is a subgraph of  $G^*$  where  $\mathcal{E}' \subseteq \mathcal{E}^*$  is the union of all maximum matchings of  $G^*$ .*

As highlighted by previous research [42, 12], if the bipartite graph is not equal to its core, then a reduction of the traceability set is possible. In the example of Figure 3, computing the core, results in the strict subset of edges  $\{(s_1, w_1), (s_2, w_2), (s_3, w_3), (s_1, w_2), (s_2, w_1)\}$ . In this case, the witness  $w_3$  has only one remaining edge and is thus completely traceable. To formalize our definition of untraceability, we therefore only consider the traceability set of a transaction after reduction of the graph. This leads us to the definition of  $k$ -untraceability for a transaction, where  $k$  is the cardinality of the traceability set.

**Definition 4 ( $k$ -untraceable transaction).** *Let  $T_w^* = (\mathcal{S}' \dot{\cup} \mathcal{N}' \dot{\cup} \{w\}, \mathcal{E})$  be a transaction such that  $\mathcal{E} = \mathcal{E}'_O \dot{\cup} \mathcal{E}'_M$  and  $T_w^* \subseteq \mathbf{core}(G^*)$ . We say that  $T_w^*$  is  $k$ -untraceable if and only if the set of ambiguous edges,  $\mathcal{E}'_O \dot{\cup} \mathcal{E}'_M$ , has cardinality  $k$  with  $k \geq 2$ .*

When an adversary  $\mathcal{A}$  observes a PDAG, the best strategy to identify the source of a transaction is to first consider the trivially traceable transactions. If no such transaction exists, the second-best strategy for  $\mathcal{A}$ , without further knowledge, is to guess uniformly at random over the edges of transactions with the smallest untraceability set.

In this view, the untraceability guarantees of the PDAG are given by the *weakest* transaction, i.e., the transactions with the least untraceability guarantees. A PDAG does not provide untraceability if any of its transaction is traceable. Otherwise the PDAG's untraceability set is tied to the smallest untraceability set of its transactions. Consequently, for a PDAG to be  $k$ -untraceable, all its transactions must at least be  $k$ -untraceable.

**Definition 5 ( $k$ -untraceable PDAG).** *A PDAG  $G$  with the corresponding bipartite graph  $G^*$  is  $k$ -untraceable if and only if all transactions  $T_W^*$  in the core of  $G^*$  are at least  $k$ -untraceable.*

### 3.3 Unlinkability

Whereas untraceability addresses the relation between input states and witnesses, unlinkability is a notion that focuses on the complete set of states in the graph. Unlinkability refers to the inability of the adversary to link two different

states or transactions together [2]. Now, if the adversary can link two states, she can also link transactions involving those states. We consequently focus on a more concise definition of unlinkability: for any two states, it is impossible to sufficiently distinguish whether they are related or not. This definition directly echoes the terminology of unlinkability of Pfizmann and Hansen [37]. We here define a *link relation* to abstract the notion of *linking*.

**Definition 6 (Link relation).** *The link relation  $L$  is an equivalence relation on the states  $\mathcal{S}$ , that is, any two states  $s_i, s_j \in \mathcal{S}$  are linked whenever  $s_i L s_j$ . We write  $[s]_L$  the equivalence class of  $\mathcal{S}$  by  $L$ , to which  $s$  belongs, i.e.,  $[s]_L := \{s_i \in \mathcal{S} : s L s_i\}$ .*

Although blockchains function differently, the main concept of unlinkability is the same for various systems. It consists of the adversary trying to link states by grouping them into those that are related. Some states may easily be linked together as they hold the same address, which typically identifies asset holders in blockchains. We therefore introduce the address relation as a refinement of the link relation.

**Definition 7 (Address relation).** *The address relation  $A$  is a refinement of the link relation  $L$  on the states  $\mathcal{S}$ , that is, any two states  $s_i, s_j \in \mathcal{S}$  hold the same address whenever  $s_i A s_j$ . We write  $[s]_A$  the equivalence class of  $\mathcal{S}$  by  $A$ , to which  $s$  belongs, i.e.,  $[s]_A := \{s_i \in \mathcal{S} : s A s_i\}$ .*

However, users are encouraged to use different addresses to break the relation between several transactions. In this case, the linking is no longer trivial, and to our knowledge, there exists no systematic way of tying states together. Research on this subject uses heuristics to associate related states together. Those heuristics are mostly patterns within or across transactions, based on some evidence of shared ownership, that lead to believe that some states may be related [19]. Similarly, we model an adversary that uses an inferred link relation  $\tilde{L}$  in her endeavor of compromising unlinkability. The adversary may use a very coarse relation that assumes that any two states are linked, or she may use finer relations such as heuristics presented in the literature or a machine-learning-based classifier. We represent here a few heuristics as an inferred linking relation  $\tilde{L}$ , and say that a subgraph  $G' \subseteq G$  is *vulnerable* to  $\tilde{L}$ , if said relation applies to  $G'$ .

The most studied heuristic is certainly the multi-input heuristic, which assumes that multiple addresses part of a single transaction are linked [33]. We express this heuristic as a linking relation:

$$\forall s_i, s_j \in \star \mathcal{E}_{CW} : s_i \tilde{L} s_j. \quad (1)$$

Another well-studied heuristic addresses the use of a change-address: when a user sends coins, she might use a newly generated address (change-address) to receive the amount surplus. If a state with a new address is part of an output set of size two, one can assume that such state is linked to at least one input state of the same transaction:

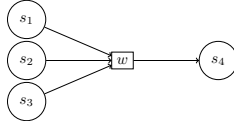
$$\exists s_i \in w \mathcal{E}_P \mathcal{S}_O, \exists s_j \in \star \mathcal{E}_{CW} : |\mathcal{S}_O| = 2 \wedge \forall s \in \mathcal{S} : s_i \notin [s]_A \implies s_i \tilde{L} s_j. \quad (2)$$

Note that this condition has different variations in the literature. The above version is stronger with  $|\mathcal{S}_O| = 2$  than with  $|\mathcal{S}_O| \geq 2$ , but the latter is usually more precise. Since the linking relations are expressed under logical statements, we can combine them to express more complex relations. For example, the multi-input heuristic, and the change-address relation can be combined to express that the state with change-address is linked to every input states:

$$\exists s_i \in w\mathcal{E}_P\mathcal{S}_O, \forall s_j \in \star\mathcal{E}_Cw : |\mathcal{S}_O| = 2 \wedge \forall s \in \mathcal{S} : s_i \notin [s]_A \implies s_i \tilde{L}s_j. \quad (3)$$

In Appendix A and in the full version of the paper [45], we translate around ten heuristics from the literature into inferred linking relations. We model here an adversary that exploits linking relations of the graph to find the different states that belong together. The equivalence class of a state  $s$  identified by the adversary is the set  $[s]_{\tilde{L}}$ , and we argue here that the dichotomy between  $[s]_L$  and  $[s]_{\tilde{L}}$  is a source of unlinkability.

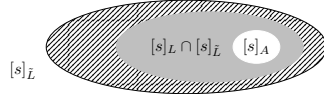
As an example, consider transaction  $T_w$  in Figure 4. We assume here that all states have a different address, that  $s_2$  and  $s_3$  are linked to  $s_4$ , and  $s_1$  is linked to neither  $s_2$ ,  $s_3$ , nor  $s_4$ . If an adversary exploits relation (1) to identify the states linked together, she will obtain  $\{s_1, s_2, s_3\}$ . In this scenario,  $s_2$  and  $s_3$  are accurately associated, but  $s_1$  is wrongfully linked to  $s_2$  and  $s_3$ . Although linking  $s_2$  and  $s_3$  together reduces unlinkability, assigning  $s_1$  to the same equivalence class has the opposite effect. The unlinkability set is therefore decreasing with any addition of a linked state into the correct equivalence class, but increasing with additions of unlinked states. In the scenario of Figure 4,  $s_1$  is singly unlinkable among three and  $s_2, s_3$  are both unlinkable among three, according to linking relation (1). The best situation for a state  $s$  in terms of unlinkability is clearly



**Fig. 4.** Example of a transaction  $T_w = (\mathcal{V}', \mathcal{E}')$  with a set of vertices  $\mathcal{V}' = \{s_1, s_2, s_3, w, s_4\}$  and a set of edges  $\mathcal{E}' = \{(s_1, w), (s_2, w), (s_3, w), (w, s_4)\}$ .

to be singly unlinkable in the wrong equivalence class, i.e.,  $s \in [s']_{\tilde{L}} \setminus [s']_L$  with  $s' \neq s$ . Conversely, the worst case scenario is when the equivalence class built by the adversary,  $[s]_{\tilde{L}}$ , accurately reflects the linking between states, i.e.,  $[s]_{\tilde{L}} \subseteq [s]_L$ .

Between those two edge cases,  $s$  still has some degree of unlinkability. This degree is directly related to the difference between the equivalence classes  $[s]_{\tilde{L}}$  and  $[s]_L$ . In Figure 5, we depict this situation, in which the outer, striped area represents the unlinked states. In contrast, the inner, gray area represents the set of states accurately linked by  $\tilde{L}$ . This set, de facto, contains the equivalent class  $[s]_A$ , since  $A$  is a refinement of  $L$  and states with the same address are trivially linkable for  $\mathcal{A}$ . Yet, the size of  $[s]_{\tilde{L}} \setminus [s]_L$  is not quite representative of



**Fig. 5.** Illustration of the unlinkability problem: the striped area represents the degree of unlinkability for state  $s$  and the gray area is the set of accurately linked states, which contains the set of states with the same address.

the degree of unlinkability of  $s$ . As an example, consider the two cases where:

1.  $|[s]_{\tilde{L}}| = 3$  and  $|[s]_{\tilde{L}} \cap [s]_L| = 2$  and
2.  $|[s]_{\tilde{L}}| = 100$  and  $|[s]_{\tilde{L}} \cap [s]_L| = 99$ .

In both cases, the cardinality of the set difference is 1. However, in the first case, an adversary has a chance of  $\binom{2}{2}/\binom{3}{2} = \frac{1}{3}$  of guessing a linked pair, whereas in the second case, this chance is increased to  $\binom{99}{2}/\binom{100}{2} \approx 0.98$ . Reciprocally, unlinkability is the ratio of pairs a state  $s$  can mix into. We define the *unlinkability score* accordingly.

**Definition 8 (Unlinkability score).** *The unlinkability score  $u$  of a state  $s$  according to an inferred linking relation  $\tilde{L}$  is the ratio of the number of unlinked pairs to all possible pairs, i.e.,  $u = \binom{|[s]_{\tilde{L}} \setminus [s]_L|}{2} / \binom{|[s]_{\tilde{L}}|}{2}$ .*

The unlinkability score  $u$  ranges from 0 to 1, where 0 denotes linkable states, and 1 represents singly-unlinkable states. We define  $u$ -unlinkability for a state with  $u > 0$ .

**Definition 9 ( $u$ -unlinkable state).** *A state  $s$  of a PDAG  $G$  is  $u$ -unlinkable according to an inferred linking relation  $\tilde{L}$  if and only if exploiting said relation on  $G$ , results in an equivalence class  $[s]_{\tilde{L}}$  with an unlinkability score  $u > 0$ .*

When observing a PDAG, the first best strategy for  $\mathcal{A}$  is to link states with the same address together, i.e.,  $\tilde{L} = A$ . In this case, for any state  $s$ ,  $[s]_{\tilde{L}}$  is subset of  $[s]_L$ , and thus the unlinkability score of all states  $s' \in [s]_{\tilde{L}}$  is 0. If  $[s]_A = \{s\}$ , the adversary must rely on non-trivial inferred linking relations  $\tilde{L}$ . These relations may nonetheless yield an unlinkability score of 0 depending on their accuracy and the topology of the PDAG. For instance, in the example of Figure 4, if  $s_1$  is linked to  $s_2$  and  $s_3$ , applying relation (1), would result in an unlinkability score of 0. Consequently, even if the second strategy is not preferable over the second, the second one may still jeopardize unlinkability completely. PDAG's unlinkability is thus tied to the smallest unlinkability score. We here define unlinkability for a PDAG according to its least unlinkable state.

**Definition 10 ( $u$ -unlinkable PDAG).** *A PDAG  $G$  is  $u$ -unlinkable according to an inferred linking relation  $\tilde{L}$  if and only if all states in  $\mathcal{S}$  are at least  $u$ -unlinkable according to  $\tilde{L}$ .*

As said before, inferred relation  $\tilde{L}$  directly refers to heuristics from the literature, and the level of unlinkability of states and ultimately, the one of the PDAG is

tied to the accuracy of this inference. However, one could imagine exploiting the different heuristics by involving unlinked states in a transaction vulnerable to  $\tilde{L}$ . This is, in a way, the principle of Coinjoin that counteracts the multi-input heuristic assumption, but we could also leverage the other heuristics and increase privacy up to their unlinkability score. In a way, the maximum unlinkability score of an inferred relation  $\tilde{L}$  can be interpreted as the *maximal extractable privacy* (MEP). In Appendix A and in the full version of the paper [45], we apply our definitions to various systems and measure their privacy guarantees. Meanwhile, we propose some extensions to our model in the following section.

## 4 Extensions

During the course of this paper, we adopted a deterministic view to describe and formalize the different privacy notions. Yet, the probabilistic standpoint might offer more flexibility and be more potent in reflecting extra information that the adversary gets from other sources (e.g., analysis from other layers, leaked databases or other attacks). Furthermore, it is probably more akin to capture machine learning-based models from the literature [32, 31, 20]. We here give a few leads to adapt our analysis to a probabilistic one by addressing each notion separately.

*Untraceability.* In Section 3, we describe the untraceability problem using a bipartite graph  $G^* = (\mathcal{S}^* \dot{\cup} \mathcal{W}^*, \mathcal{E}^*)$ . As part of the extension to a probabilistic model, we could assign weights to  $\mathcal{E}^*$ . Let  $\omega : \mathcal{E}^* \rightarrow [0, 1] \subseteq \mathbb{R}$  be a map that assigns a real number between 0 and 1 to each edge. The weight  $\omega(e)$  of edge  $e = (s, w)$  represents the probability that the underlying state  $s$  is consumed in the transaction involving witness  $w$ . This means that the sum of weights for edges of a single witness sums up to 1. The adversary is therefore able to assign a probability distribution based on some information (e.g., time or value-based correlations) or an inference of a machine learning model.  $\mathcal{A}$  can then output a guess according to the weights of the edges.

*Unlinkability.* Although the inferred linking relation  $\tilde{L}$  can already abstract a probabilistic classifier, we may go further by modifying this relation into a fuzzy relation  $\mu$  such that  $\mu : \mathcal{S} \times \mathcal{S} \rightarrow [0, 1] \subseteq \mathbb{R}$ . This relation gives a membership indicator between two states  $s_i$  and  $s_j$ , where 0 means that  $s_i$  and  $s_j$  are not related whereas 1 means they are. This would allow one state  $s$  to belong to two distinct clusters or to measure uncertainty in the clustering output.

## 5 Related Work

Regarding untraceability, Möser et al. [32] conduct an empirical analysis in the Monero blockchain. They identify two phenomena that jeopardize the untraceability guarantees of Monero. Möser [32] notice that 0-mixin transactions not only imply full traceability of such users, but also pose a privacy risk for other

users that include poorly mixed outputs as mixins in subsequent transactions. Fortunately, since October 2018, the mixins input number is set globally to eleven [29], consequently alleviating this issue. Yet, Möser et al. [32] make a second observation regarding the sampling distribution of mixins inputs. The distribution of the sampling algorithm to choose the different mixins does not match the distribution of the coins spent in the transaction. The sampling distribution is uniform over the set of available transaction outputs, irrespective of the age of coins, whereas new coins are more likely to be spent directly. Old mixin coins can thus be considered decoys, reducing the true spender’s anonymity set. Following this study, the Monero community also addressed this issue with a better sampling algorithm [13].

As for unlinkability, various works devise empirical behavioral patterns in transactions – called heuristics – to link addresses together. Nakamoto [33] describes the most widely studied heuristic in the original Bitcoin white paper. This heuristic allows linking to the same entity, multiple addresses used as input to the same transaction [33]. Studies find it to be relatively accurate [2] and able to reveal hidden clusters of addresses [28]. The Coinjoin mechanism mitigates the threat posed by this heuristic by having different users join their coins into a single transaction to blur the links between several of their addresses. Yet, Goldfeder et al. [15] devise a counter-heuristic to detect such a technique by considering multi-input-multi-output transactions as part of a Coinjoin mixing. Victor [41] devises heuristics for specific patterns in the Ethereum and Wang et al. [43] for cases where users utilize mixers to break links between their transactions. Finally, Kappos et al. [19] expose heuristics specifically for Zcash. Our work is also close to any literature review of deanonymization heuristics such as the systematization of knowledge of [9]. Furthermore, Meiklejohn and Orlandi [27] devise a metric similar to the unlinkability score for Coinjoin. We yet believe our framework and metric to be more generally applicable to any scheme and adversarial classifier.

Concerning graph-based analysis of privacy, Ober et al. [35] conduct an empirical study of crucial properties of the Bitcoin transaction graph, such as unlinkability and coin dormancy. They use the multi-input heuristic mentioned above to merge addresses that belong to the same entity. As a result, they discover an impressive amount of public keys that belong together while only considering a small subset of the full blockchain history. Androulaki et al. [2] carry out a similar study through a simulator that mimics the user of Bitcoin within a university. In this setting, the results of the heuristics can be compared to the ground truth. It turns out that almost 40% of user profiles can be recovered even though recommended privacy measures are applied. Also, Atzei et al. [3] develop a model to formally prove the fulfillment of some properties in the Bitcoin blockchain [3]. Conversely, Amarasinghe et al. [1] employ a common, universal framework to characterize the various aspects of anonymity provided by different blockchain implementations [1].

Finally, our work is also related to research on privacy notions in anonymous communication networks. Kuhn et al. [23] provide an in-depth analysis on

the matter [23]. They study privacy notions in anonymous communication and their relations and devise a formal hierarchy to classify each of these properties. Moreover, Henry et al. [16] aim at formalizing anonymous blacklisting systems.

## 6 Conclusion

This work has tackled the definition of two main blockchain privacy notions: untraceability and unlinkability. To do so, we first extended the TDAG [7] to capture privacy-preserving blockchains (PDAG) and, secondly, gave consistent definitions to these notions according to PDAG’s elements. This allowed us to model and compare blockchain implementations and unify literature results in Appendix A. During the course of this work, we also have discovered that the PDAG gives intuitive definitions to each notion and proposes a way to reason about them.

**Acknowledgment.** We thank the anonymous reviewers for helpful suggestions and feedback. Special thanks go to Jayamine Alupotha and Mathieu Gestin for interesting discussions about blockchain privacy and valuable comments. We are also grateful to David Lehnher for fruitful advice and comments. Duc V. Le has been supported by a grant from Protocol Labs to the University of Bern.

## References

1. Amarasinghe, N., Boyen, X., McKague, M.: The complex shape of anonymity in cryptocurrencies: Case studies from a systematic approach. In: *Financial Cryptography* (1). *Lecture Notes in Computer Science*, vol. 12674, pp. 205–225. Springer (2021)
2. Androulaki, E., Karame, G., Roeschlin, M., Scherer, T., Capkun, S.: Evaluating user privacy in Bitcoin. In: *Financial Cryptography*. *Lecture Notes in Computer Science*, vol. 7859, pp. 34–51. Springer (2013)
3. Atzei, N., Bartoletti, M., Lande, S., Zunino, R.: A formal model of Bitcoin transactions. In: *Financial Cryptography*. *Lecture Notes in Computer Science*, vol. 10957, pp. 541–560. Springer (2018)
4. Ben-Sasson, E., Chiesa, A., Garman, C., Green, M., Miers, I., Tromer, E., Virza, M.: Zerocash: Decentralized anonymous payments from Bitcoin. In: *IEEE Symposium on Security and Privacy*. pp. 459–474. IEEE Computer Society (2014)
5. Biryukov, A., Khovratovich, D., Pustogarov, I.: Deanonymisation of clients in bitcoin P2P network. In: *CCS*. pp. 15–29. ACM (2014)
6. Biryukov, A., Pustogarov, I.: Bitcoin over tor isn’t a good idea. In: *IEEE Symposium on Security and Privacy*. pp. 122–134. IEEE Computer Society (2015)
7. Cachin, C., De Caro, A., Moreno-Sanchez, P., Tackmann, B., Vukolić, M.: The transaction graph for modeling blockchain semantics. *Cryptoeconomic Systems* **1**(1) (2021), <https://doi.org/10.21428/58320208.a12c57e6>, preliminary version appears as *Cryptology ePrint Archive*, Report 2017/1070, 2017
8. Chaum, D.: The dining cryptographers problem: Unconditional sender and recipient untraceability. *J. Cryptol.* **1**(1), 65–75 (1988)



9. Deuber, D., Ronge, V., Rückert, C.: Sok: Assumptions underlying cryptocurrency deanonymizations. *Proc. Priv. Enhancing Technol.* **2022**(3), 670–691 (2022)
10. Dulmage, A.L., Mendelsohn, N.S.: Coverings of bipartite graphs. *Canadian Journal of Mathematics* **10**, 517–534 (1958). <https://doi.org/10.4153/CJM-1958-052-0>
11. Edman, M., Sivrikaya, F., Yener, B.: A combinatorial approach to measuring anonymity. In: ISI. pp. 356–363. IEEE (2007)
12. Egger, C., Lai, R.W.F., Ronge, V., Woo, I.K.Y., Yin, H.H.F.: On defeating graph analysis of anonymous transactions. *Proc. Priv. Enhancing Technol.* **2022**(3), 538–557 (2022)
13. Ehrenhofer, J.: Response to "an empirical analysis of traceability in the Monero blockchain", Version 2. <https://www.getmonero.org/2018/03/29/response-to-an-empirical-analysis-of-traceability.html> (2018)
14. Gierlichs, B., Troncoso, C., Díaz, C., Preneel, B., Verbaauwhede, I.: Revisiting a combinatorial approach toward measuring anonymity. In: WPES. pp. 111–116. ACM (2008)
15. Goldfeder, S., Kalodner, H.A., Reisman, D., Narayanan, A.: When the cookie meets the blockchain: Privacy risks of web payments via cryptocurrencies. *Proc. Priv. Enhancing Technol.* **2018**(4), 179–199 (2018)
16. Henry, R., Goldberg, I.: Formalizing anonymous blacklisting systems. In: IEEE Symposium on Security and Privacy. pp. 81–95. IEEE Computer Society (2011)
17. Hopwood, D., Bow, S., Hornby, T., Wilcox, N.: Zcash protocol specification. <https://zips.z.cash/protocol/protocol.pdf> (2021)
18. Kalodner, H.A., Möser, M., Lee, K., Goldfeder, S., Plattner, M., Chator, A., Narayanan, A.: Blocksci: Design and applications of a blockchain analysis platform. In: USENIX Security Symposium. pp. 2721–2738. USENIX Association (2020)
19. Kappos, G., Yousaf, H., Maller, M., Meiklejohn, S.: An empirical analysis of anonymity in Zcash. In: USENIX Security Symposium. pp. 463–477. USENIX Association (2018)
20. Kappos, G., Yousaf, H., Stütz, R., Rollet, S., Haslhofer, B., Meiklejohn, S.: How to peel a million: Validating and expanding Bitcoin clusters. In: USENIX Security Symposium. pp. 2207–2223. USENIX Association (2022)
21. Koe, Alonso, K.M., Noether, S.: Zero to Monero. Online, second edn. (2020), <https://web.getmonero.org/library/Zero-to-Monero-2-0-0.pdf>
22. Koshy, P., Koshy, D., McDaniel, P.D.: An analysis of anonymity in bitcoin using P2P network traffic. In: Financial Cryptography. Lecture Notes in Computer Science, vol. 8437, pp. 469–485. Springer (2014)
23. Kuhn, C., Beck, M., Schiffner, S., Jorswieck, E.A., Strufe, T.: On privacy notions in anonymous communication. *Proc. Priv. Enhancing Technol.* **2019**(2), 105–125 (2019)
24. Kumar, A., Fischer, C., Tople, S., Saxena, P.: A traceability analysis of Monero’s blockchain. In: ESORICS (2). Lecture Notes in Computer Science, vol. 10493, pp. 153–173. Springer (2017)
25. Le, D.V., Gervais, A.: AMR: autonomous coin mixer with privacy preserving reward distribution. In: AFT. pp. 142–155. ACM (2021)
26. Maxwell, G.: CoinJoin: Bitcoin privacy for the real world. <https://bitcointalk.org/index.php?topic=279249.0> (2013)
27. Meiklejohn, S., Orlandi, C.: Privacy-enhancing overlays in bitcoin. In: Financial Cryptography Workshops. Lecture Notes in Computer Science, vol. 8976, pp. 127–141. Springer (2015)

28. Meiklejohn, S., Pomarole, M., Jordan, G., Levchenko, K., McCoy, D., Voelker, G.M., Savage, S.: A fistful of Bitcoins: Characterizing payments among men with no names. *login Usenix Mag.* **38**(6) (2013)
29. Monero Project: Monero 0.13.0 "Beryllium Bullet" Release. <https://www.getmonero.org//2018/10/11/monero-0.13.0-released.html> (2018)
30. Monero Project: Monero v.0.18.0 release. <https://github.com/monero-project/monero/releases/tag/v0.18.0.0> (2023)
31. Möser, M., Narayanan, A.: Resurrecting address clustering in Bitcoin. In: *Financial Cryptography. Lecture Notes in Computer Science*, vol. 13411, pp. 386–403. Springer (2022)
32. Möser, M., Soska, K., Heilman, E., Lee, K., Heffan, H., Srivastava, S., Hogan, K., Hennessey, J., Miller, A., Narayanan, A., Christin, N.: An empirical analysis of traceability in the monero blockchain. *Proc. Priv. Enhancing Technol.* **2018**(3), 143–163 (2018)
33. Nakamoto, S.: Bitcoin: A peer-to-peer electronic cash system. <https://bitcoin.org/bitcoin.pdf> (2008)
34. Noether, S., Mackenzie, A., Research Lab, T.M.: Ring Confidential Transactions. *ledger* **1**, 1–18 (2016). <https://doi.org/10.5195/ledger.2016.34>
35. Ober, M., Katzenbeisser, S., Hamacher, K.: Structure and anonymity of the Bitcoin transaction graph. *Future Internet* **5**(2), 237–250 (2013)
36. Pertsev, A., Semenov, R., Storm, R.: Tornado cash privacy solution (2019)
37. Pfitzmann, A., Köhntopp, M.: Anonymity, unobservability, and pseudonymity - A proposal for terminology. In: *Workshop on Design Issues in Anonymity and Unobservability. Lecture Notes in Computer Science*, vol. 2009, pp. 1–9. Springer (2000)
38. Rivest, R.L., Shamir, A., Tauman, Y.: How to leak a secret. In: *ASIACRYPT. Lecture Notes in Computer Science*, vol. 2248, pp. 552–565. Springer (2001)
39. van Saberhagen, N.: CryptoNote. <https://bytecoin.org/old/whitepaper.pdf> (2013)
40. Venkatakrishnan, S.B., Fanti, G., Viswanath, P.: Dandelion: Redesigning the bitcoin network for anonymity. *Proc. ACM Meas. Anal. Comput. Syst.* **1**(1), 22:1–22:34 (2017)
41. Victor, F.: Address clustering heuristics for Ethereum. In: *Financial Cryptography. Lecture Notes in Computer Science*, vol. 12059, pp. 617–633. Springer (2020)
42. Vijayakumaran, S.: Analysis of cryptonote transaction graphs using the dulmage-mendelsohn decomposition. In: *AFT. LIPIcs*, vol. 282, pp. 28:1–28:22. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2023)
43. Wang, Z., Chaliasos, S., Qin, K., Zhou, L., Gao, L., Berrang, P., Livshits, B., Gervais, A.: On how zero-knowledge proof blockchain mixers improve, and worsen user privacy. In: *WWW*. pp. 2022–2032. ACM (2023)
44. Wang, Z., Cirkovic, M., Le, D.V., Knottenbelt, W.J., Cachin, C.: Pay less for your privacy: Towards cost-effective on-chain mixers. In: *AFT. LIPIcs*, vol. 282, pp. 16:1–16:25. Schloss Dagstuhl - Leibniz-Zentrum für Informatik (2023)
45. Wicht, F.X., Wang, Z., Le, D.V., Cachin, C.: A transaction-level model for blockchain privacy. *Cryptology ePrint Archive*, Paper 2023/1902 (2023), <https://eprint.iacr.org/2023/1902>, <https://eprint.iacr.org/2023/1902>
46. Yu, Z., Au, M.H., Yu, J., Yang, R., Xu, Q., Lau, W.F.: New empirical traceability analysis of CryptoNote-style blockchains. In: *Financial Cryptography. Lecture Notes in Computer Science*, vol. 11598, pp. 133–149. Springer (2019)

## A Applications

In this section, we analyze four existing add-on privacy solutions and privacy-preserving blockchains using the privacy notions in the PDAG model.

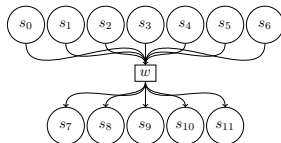
### A.1 Add-on Privacy Solutions

The two most prominent blockchains, Bitcoin and Ethereum, do not provide inherent privacy features. Users are therefore encouraged to regularly change their address to break the link between several of their transactions. However, in this case, users' wealth is scattered among multiple addresses. If users want to conduct transactions that exceed the amount held at one address, they need to merge multiple of their addresses into a single transaction. Yet, as discussed in the unlinkability subsection, this pattern is noticeable. Merging transactions are easily identified and input states are deemed to be related. We here present two add-on privacy solutions that mitigate this issue: Coinjoin and Tornado Cash. These two mechanisms allows for breaking the link between several transactions. Yet, transactions using these techniques still reveal the state consumption. Hence, according to our definitions, these two privacy-enhancing solutions do not provide untraceability, but only unlinkability.

**Coinjoin Background.** Coinjoin [26] is a mixing protocol that allows multiple users to merge their coins into a single transaction. In order to enhance unlinkability,  $U$  users agree on a standardized output size and collectively contribute inputs totaling at least that size. The resulting transaction would consist of  $U$  outputs, each of the agreed-upon size, and potentially  $U$  additional change outputs if some users provided inputs exceeding the target. All participants would sign the transaction, and the transaction can be conducted. *In fine*, the goal of this protocol is to escape the multi-input heuristic presented in subsection 3.3.

*Execution.* Consider the transaction  $T_w$  depicted on Figure 6. This transaction transfers funds of three different users:  $u_1$ ,  $u_2$ , and  $u_3$ . User  $u_1$  sends coins from her two addresses,  $s_0, s_1$ , to  $s_7$ , and gets the rest amount in a change-address,  $s_8$ . Analogously, the second user,  $u_2$ , merges two addresses,  $s_2$  and  $s_3$ , to a single address,  $s_9$ . Finally,  $u_3$  transfers funds from three addresses,  $s_4, s_5$ , and  $s_6$ , to address  $s_{10}$ , and receives the surplus to  $s_{11}$ . The witness  $w$  holds the verification for the seven input signatures. In this case, if  $\mathcal{A}$  uses the multi-input heuristic, she will get an equivalence class of  $\{s_1, s_2, s_3, s_4, s_5, s_6\}$ , and the unlinkability scores will be  $u_1 = u_2 = u_3 = u_4 = \frac{10}{21}$  and  $u_5 = u_6 = \frac{6}{21}$ .

*Privacy guarantees.* According to our definitions, Coinjoin does not provide untraceability guarantees. However, a single Coinjoin transaction under the multi-input heuristic provides unlinkability within the bound of the number of inputs and users' contributions.



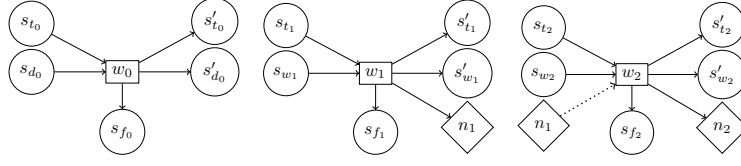
**Fig. 6.** An illustration of a Coinjoin transaction,  $T_w$ , that merges funds from seven addresses ( $s_0$  to  $s_6$ ) to five different ones ( $s_7$  to  $s_{11}$ ).

**Lemma 1 (Coinjoin unlinkability).** *A Coinjoin transaction  $T_w = \mathcal{S}_I w \mathcal{S}_O$  containing at least two disjoint equivalence classes by  $L$ , from which  $U$  is the largest and  $u$  the smallest, provides unlinkability score within the range of  $[(\binom{|S_I|-|U|}{2} / \binom{|S_I|}{2}), (\binom{|S_I|-|u|}{2} / \binom{|S_I|}{2})]$  under the multi-input relation  $\tilde{L}$ .*

*Proof.* The equivalence class by the multi-input relation  $\tilde{L}$  is  $[s]_{\tilde{L}} = \mathcal{S}_I$ , for all  $s \in \mathcal{S}_I$ . We have at least two disjoint equivalence classes  $[s_1]_{\tilde{L}} \subseteq [s_1]_{\tilde{L}}$  and  $[s_2]_{\tilde{L}} \subseteq [s_2]_{\tilde{L}}$ . Without loss of generality, we assume that  $|[s_1]_{\tilde{L}}| \leq |[s_i]_{\tilde{L}}|$  and that  $|[s_2]_{\tilde{L}}| \geq |[s_i]_{\tilde{L}}|$  for all  $s_i \in \mathcal{S}_I$ . The score  $u_1 = (\binom{|S_I|-|U|}{2} / \binom{|S_I|}{2})$  is therefore the smallest since  $U$  is the largest equivalence class. Additionally, the score  $u_2 = (\binom{|S_I|-|u|}{2} / \binom{|S_I|}{2})$  is the largest since  $u$  is the smallest equivalence class. Consequently, the unlinkability score lies within the range  $[u_1, u_2]$ .

**Tornado Cash Background.** Tornado Cash is a zero-knowledge proof (ZKP) mixer [44, 25, 36] that runs on smart-contract-enabled blockchains such as Ethereum. A ZKP-mixer is a smart contract that supports the mixing service of a fixed amount of coins. Users deposit coins in the mixer and receive a deposit note. At a later point, users can withdraw their coins by proving knowledge of the deposit note. The withdrawal is done in a zero-knowledge fashion to conceal any link between depositors and withdrawers.

*Execution.* Figure 7 shows three transactions that involve Tornado Cash:  $T_{w_0}$ ,  $T_{w_1}$  and  $T_{w_2}$ . We assume a topological order of  $w_0 \prec w_1 \prec w_2$ . First,  $T_{w_0}$  represents the deposit of coins from an address ( $s_{d_0}$ ). The smart contract ( $s_{t_0}$ ) records the deposit note as a new entry ( $s'_{t_0}$ ). The depositor also needs to pay a fee ( $s_{f_0}$ ) for invoking the smart contract. The witness  $w_0$  holds the verification of the signature for the smart contract invocation. In this example, we assume that multiple deposits analogous to  $t_0(w_0)$  take place. The second transaction  $T_{w_1}$  represents an address ( $s_{w_1}$ ) that invokes the smart contract ( $s_{t_1}$ ) to withdraw *some* previous deposit ( $s'_{w_1}$ ). Again, this operation incurs a fee ( $s_{f_1}$ ). Also, this transaction generates a withdrawal commitment ( $n_1$ ), preventing the withdrawal of the same deposit. Finally, the third transaction  $T_{w_2}$  represents the first withdrawal that comes after withdrawal  $T_{w_1}$ . In this case, an address ( $s_{w_2}$ ) executes the smart contract ( $s_{t_2}$ ) to withdraw a deposit ( $t_2(w_2)$ ). A fee ( $s_{f_2}$ ) is collected. The previously generated withdrawal commitment ( $n_1$ ) is read by the smart contract for any conflict. The updated list of commitments ( $n_2$ ) is then generated.



**Fig. 7.** Illustration of three transactions  $T_{w_0}$ ,  $T_{w_1}$ ,  $T_{w_2}$  involving Tornado Cash. The leftmost graph,  $T_{w_0}$ , represents a deposit. The center graph,  $T_{w_1}$ , is an example of a first withdrawal, while the rightmost depicts a subsequent withdrawal to  $T_{w_1}$ .

*Privacy guarantees.* Similarly to Coinjoin, Tornado Cash does not provide untraceability, but unlinkability. However, the ZKP-mixer does not consist of a single transaction, but of multiple, seemingly unrelated, deposits and withdrawals. Wang et al. [43] devise several heuristics as a way to link the different states involved in deposits and withdrawals. The first heuristic they describe is the *address reuse*, when both depositor and withdrawer use the same address. In this case, the adversary uses the refinement of the linking relation  $A$  to link depositors to withdrawers. This results in an unlinkability score of 0, since  $[s]_A \subseteq [s]_L$ , for all  $s \in \mathcal{S}$ .

A second heuristic, Wang et al. describe, is the reuse of the deposit address to pay the withdrawal fees. Let  $F$  be the transactional relation between a state  $s$  and a witness  $w$ , such that,  $s$  pays the fee of transaction involving  $w$ , whenever  $sFw$ . We express this second heuristic according to an inferred linking relation  $\tilde{L}$ , and denote  $w_d$  the witness involved in the deposit and  $w_w$  in the withdrawal:

$$\forall s_i \in \star \mathcal{E}_C w_d, s_j, s_f \in \star \mathcal{E}_C w_w : s_f F w_w \wedge s_f A s_i \implies s_j \tilde{L} s_f \wedge s_f L s_i. \quad (4)$$

This heuristic identifies three different states among which two with the same address are linked.

**Lemma 2.** *Any state involved in a Tornado Cash transaction  $T_w$  vulnerable to inferred relation  $\tilde{L}$  in (4) is either linkable or has an unlinkability score of  $\frac{2}{3}$  for state  $s \notin [s']_A$  for all  $s' \in T_w, s' \neq s$ .*

*Proof.* Say we have three states involved in transaction  $T_w$ :  $s_d$ ,  $s_f$ , and  $s_w$ . Transaction  $T_w$  is vulnerable to inferred relation  $\tilde{L}$  described in (4), and thus  $s_f A s_d$ . If  $s_w \in [s_d]_L$ , then all states are linkable. Otherwise,  $s_w$  has an unlinkability score of  $(\binom{3}{2} - \binom{2}{2}) / \binom{3}{2} = \frac{2}{3}$ .

A third heuristic involving Tornado Cash is when the withdrawal and the deposit addresses are involved in a subsequent transaction. In this case, one can assume that both addresses are linked:

$$\forall s_i \in \star \mathcal{E}_C w_d, s_j \in \star \mathcal{E}_C w_w \subseteq T_w : s_i, s_j \in T_l \succ T_w \implies s_i \tilde{L} s_j. \quad (5)$$

**Lemma 3.** *Any state involved in a Tornado Cash transaction  $T_w$  vulnerable to inferred relation  $\tilde{L}$  in (5) is either linkable or has an unlinkability score of 1.*

*Proof.* Take the two states  $s_i$  and  $s_j$  involved transaction  $T_w$ . If  $s_i \in [s_j]_L$ , then all states are linkable. Otherwise, both states have an unlinkability score of  $\binom{2}{2}/\binom{2}{2} = 1$ .

Wang et al. [43] give a fourth heuristic where an address splits its funds into three different addresses, before them depositing into the mixer. This pattern leads to inferring a link between all four addresses:

$$\forall s_i \in \star\mathcal{E}_C\mathcal{S}_O, s_j \in \star\mathcal{E}_Cw_d : |\mathcal{S}_O| = 3 \wedge s_j \in \mathcal{S}_O \implies s_i\tilde{L}s_j. \quad (6)$$

**Lemma 4.** *Any state  $s$  involved in a Tornado Cash transaction  $T_w$  vulnerable to inferred relation  $\tilde{L}$  in (6) has an unlinkability score of  $u = \frac{12-l(l-1)}{12}$ , with  $l = |[s]_L \cap [s]_{\tilde{L}}|$ .*

*Proof.* Say we have four identified states:  $s_1, s_2, s_3$ , and  $s_4$ . The equivalence class by  $\tilde{L}$  over  $T_w$  is  $[s]_{\tilde{L}} = \{s_1, s_2, s_3, s_4\}$ , for all  $s \in T_w$ . We thus have a score  $u = \left(\binom{4}{2} - \binom{l}{2}\right)/\binom{4}{2} = \frac{12-l(l-1)}{12}$ , with  $l = |[s]_L \cap [s]_{\tilde{L}}|$ .

Finally, we can also infer a link between two addresses if one deposits into the same set of mixers as the other withdraws from [43]:

$$\forall s_i \in \star\mathcal{E}_Cw_{d_1}, s_j \in \star\mathcal{E}_Cw_{d_2}, s_k \in \star\mathcal{E}_Cw_{w_1}, s_l \in \star\mathcal{E}_Cw_{w_2} : \quad (7)$$

$$s_iAs_j \wedge s_kAs_l \implies s_i\tilde{L}s_k.$$

**Lemma 5.** *Any state  $s$  involved in a Tornado Cash transaction  $T_w$  vulnerable to inferred relation  $\tilde{L}$  in (7) is either linkable or has an unlinkability score of  $u = \left(\binom{2l}{2} - \binom{l}{2}\right)/\binom{2l}{2}$ , with  $l = |[s]_L \cap [s]_{\tilde{L}}|$ .*

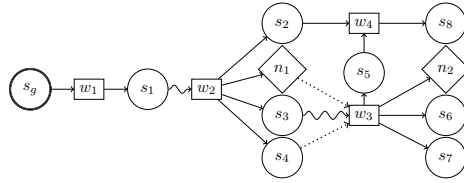
*Proof.* Say we have two equivalence classes: the deposit states  $[s_1]_A$  and the withdrawer states  $[s_2]_A$ . If  $[s_1]_A \subseteq [s]_L$  and  $[s_2]_A \subseteq [s]_L$ , for all states  $s \in T_w$ , then all states are linkable. Otherwise, we have two disjoint equivalence classes by  $L$  in  $[s]_{\tilde{L}}$  of size  $l$ , and the score for each  $s \in T_w$  is  $u = \left(\binom{2l}{2} - \binom{l}{2}\right)/\binom{2l}{2}$ .

## A.2 Monero

**Background.** Monero is a privacy-preserving blockchain that follows the UTXO model. Much like Bitcoin, each transaction transfers some inputs into some outputs. Each output is assigned to a public key and, to spend it, one must sign a message with the corresponding private key, thus proving ownership of the output. While Bitcoin discloses the source of transactions, Monero hides this information with a ring signature [38]. This mechanism is composed of a signature and a set of public keys. The signature is valid if it was produced using the private key of one of the corresponding public keys, without revealing which one of them. Since originators of transactions are hidden, additional components must also be implemented to prevent double-spending. In Monero, for each output spent, one has to produce a key image with the private key of the output. Any

attempt to spend the same output twice will produce the same key image, exposing the attempted double-spending. In addition, Monero generates one-time addresses for each output. Any output in the Monero blockchain has therefore a distinct address. Finally, since 2017, Monero has further strengthened its privacy properties by introducing ring confidential transactions (RingCTs) [34] to ensure that the transacted amount is not disclosed.

*Sample execution.* Figure 8 shows the PDAG of a Monero execution with four transactions. First,  $T_{w_1}$  represents the transfer of the mining reward to a user  $u$  that successfully mined the first block ( $s_1$ ). Second,  $T_{w_2}$  is a transaction in which  $u$  transfers some of her coins ( $s_1$ ) to another address ( $s_4$ ). The transaction fee is modeled as another state ( $s_2$ ) and the remainder of  $u$ 's coins are transferred to another address ( $s_3$ ). In that case,  $w_2$  represents the validation of the signature of  $s_1$ , and the key images produced are encapsulated in a nullifier ( $n_1$ ). The third transaction  $T_{w_3}$  represents the transfer of coins held by an address ( $s_3$ ) to two different ones ( $s_6, s_7$ ). Again, a fee ( $s_5$ ) is incurred. The witness  $w_3$  represents the verification of the signature from  $s_3$  and of the conflict-freedom with the key image in  $n_1$ . A nullifier,  $n_2$ , consists of the previous key image as well as the newly produced one. Finally, the mining transaction  $T_{w_4}$  takes place with the transfer of the fees ( $s_2, s_5$ ) and the mining reward to the miner ( $s_8$ ).



**Fig. 8.** Illustration of a PDAG  $G$  that depicts a Monero execution of four transactions:  $T_{w_1}$ ,  $T_{w_2}$ ,  $T_{w_3}$  and  $T_{w_4}$ . Each transaction  $T_i$  is the closed neighborhood subgraph  $N_G[w_i] \subset G$ .

*Privacy guarantees.* As highlighted above, Monero has built-in untraceability mechanism with ring signature. Since version 0.18.0.0 *Fluorine Fermi* of the Monero protocol, the ring size is set to 16 [30].

**Lemma 6.** *Consider a PDAG  $G$  modeling a Monero execution since version 0.18.0.0. If the corresponding bipartite graph  $G^*$  is equal to its core, then  $G$  is 16-untraceable.*

*Proof.* By contradiction, assume that  $G$  is  $k$ -untraceable with  $k < 16$ . That means that there exists a witness  $w$  in  $G^*$ , with  $k$  ambiguous edges. Since the ring signature of Monero protocol version 0.18.0.0 uses 16 inputs, there exists a graph decomposition that reduces the number of edges of  $w$  to  $k$ . Yet, that contradicts the assumption that  $G^* = \mathbf{core}(G^*)$ .

Regarding unlinkability, Monero uses one-time addresses so that any two states in the PDAG are not part of the same equivalence class by  $A$ . Monero has yet been vulnerable to multiple heuristics over the years, many of which have become obsolete thanks to incremental updates of the core protocol. As far as our knowledge goes, the multi-output heuristic formulated by Kumar et al. [24] remains the only one relevant at the time of writing. This heuristic relies on the assumption that when two outputs from a single transaction are observed together in distinct input rings of another transaction, it indicates that these two outputs are related. Let  $R$  be an equivalence relation on the states  $\mathcal{S}$  that denotes whether two states  $s_i, s_j \in \mathcal{S}$  are part of the same ring, whenever  $s_i R s_j$ . We write  $[s]_R$  the equivalence class of  $\mathcal{S}$  by  $R$ , to which  $s$  belongs. We express this heuristic as an inferred linking relation  $\tilde{L}$ :

$$\forall s_i, s_j \in w_1 \mathcal{E}_P \mathcal{S}_O : s_i, s_j \in \mathcal{S}_I \mathcal{E}_O w_2 \wedge s_i \notin [s_j]_R \implies s_i \tilde{L} s_j \quad (8)$$

**Lemma 7.** *Any state  $s$  involved in a subgraph  $G'$  vulnerable to inferred relation  $\tilde{L}$  in (8) has an unlinkability score of  $(\binom{(\mathcal{S}_O \cap \mathcal{S}_I) \setminus [s]_R \setminus [s]_L}{2}) / \binom{(\mathcal{S}_O \cap \mathcal{S}_I) \setminus [s]_R}{2}$ .*

*Proof.* The inferred relation  $\tilde{L}$  in (8) gives us an equivalence class  $[s]_{\tilde{L}}$  of all the output states  $\mathcal{S}$  taken as input to a subsequent transaction  $\mathcal{S}_I$  while omitting the ones in the same ring, thus  $(\mathcal{S}_I \cap \mathcal{S}_O) \setminus [s]_R$ . The unlinkability score is thus as follows:  $(\binom{(\mathcal{S}_O \cap \mathcal{S}_I) \setminus [s]_R \setminus [s]_L}{2}) / \binom{(\mathcal{S}_O \cap \mathcal{S}_I) \setminus [s]_R}{2}$ .

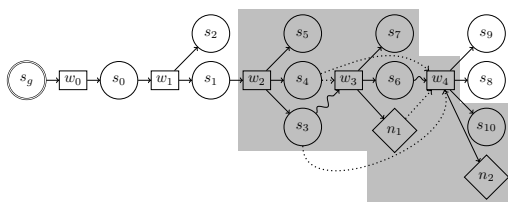
### A.3 Zcash

**Background.** Zcash behaves similarly to Bitcoin, but with the added feature of providing enhanced privacy for its users. To address this, Zcash implements a shielded pool, which is a blockchain subset in which senders, recipients, and the transferred amounts are all hidden. ZK-SNARK is used by the coin sender to prove knowledge of the secret key of this coin, without revealing which coin is actually spent. The source of a transaction is therefore obfuscated within the whole set of shielded coins. Furthermore, a valid transaction ensures that coins cannot be spent twice by generating nullifiers. The sender also conceals the recipient by using *key-private* asymmetric encryption on the recipient address. In turn, the recipient scans the blockchain for coins that belong to her. Finally, since Zcash possesses both a transparent pool, in which addresses begin with  $t$ , and a shielded pool with privacy features, in which addresses are prefixed with  $z$ , the blockchain supports four different types of transactions:  $t$ -to- $t$ ,  $t$ -to- $z$ ,  $z$ -to- $z$  and  $z$ -to- $t$ .

*Sample execution.* In Figure 9, we show five transactions part of a Zcash execution. The gray area highlights the shielded pool. The first transaction  $T_{w_0}$  is the pre-mine, minting all available coins into an address containing unmined coins ( $s_0$ ). The second transaction  $T_{w_1}$  represents the transfer of some unmined coins ( $s_0$ ) to the Zcash address of a user  $u$  who successfully mined the first block ( $s_1$ ). The remaining unmined coins are left in a distinct state ( $s_2$ ). Subsequent



transactions take place in the shielded pool. The third transaction  $T_{w_2}$  has user  $u$  send her coins to two shielded addresses  $(s_3, s_4)$ . The witness  $w_2$  represents the verification of the signature of the spender  $(s_1)$ . A fee  $(s_5)$  is also incurred in that case. The next transaction  $T_{w_3}$  transfers coins from a shielded address  $(s_3)$  to another one  $(s_6)$ , minus a fee  $(s_7)$ . The witness  $w_3$  holds the verification for the zero-knowledge proof. All the shielded addresses are part of this transaction  $(s_3, s_4)$  since no information is leaked on the sender. Furthermore, a nullifier  $(n_1)$  is produced by the transaction to prevent double-spending. The fifth transaction  $T_{w_4}$  represents the transfer from one shielded address  $(s_6)$  to another shielded address  $(s_{10})$  and to two transparent one  $(s_8, s_{10})$ . A fee  $(s_9)$  is again collected. The verification of the zero-knowledge proof and the conflict-freedom takes place at the witness  $w_4$ . It takes the previous nullifier  $n_1$  and produces a new one  $n_2$ ; the first being included in the latter.



**Fig. 9.** Illustration of a PDAG  $G$  that depicts a Zcash execution of five transactions:  $T_{w_0}, T_{w_1}, T_{w_2}, T_{w_3}$  and  $T_{w_4}$ . Each transaction  $T_i$  is the closed neighborhood subgraph  $N_G[w_i] \subset G$ . The gray area represents the shielded pool of the Zcash blockchain.

*Privacy guarantees.* Zcash’s transparent pool does not provide any untraceability guarantees. However, the shielded pool has inherent mechanisms that promote untraceability. In contrast to Monero, the untraceability set of each transaction is implicit, and it consists of any previously produced states.

**Lemma 8.** *If the PDAG  $G$  modeling a Zcash execution of the shielded pool with the corresponding bipartite graph  $G^*$  is equal to its core, then  $G$  is  $p$ -untraceable, with  $p = \min(|\{s \mid s \prec w, \forall w \in \mathcal{W}\}|)$ .*

*Proof.* By contradiction, assume that  $G$  is  $k$ -untraceable, with  $k < p$ . That means that there exists a transaction  $T_w$ , with  $k < p$  incoming ambiguous edges, i.e., with strictly fewer states taken as inputs than preceding states. However, Zcash shielded pool does not leak any information on the sender, and the possible senders are therefore the set of previously produced states. That means, that there exists a graph decomposition that reduces the number of edges of  $w$  to  $k$ . Yet, that contradicts the assumption that  $G^* = \mathbf{core}(G^*)$ .

Regarding unlinkability, Zcash’s shielded pool suffers mainly from the interactions with the transparent pool. Kappos et al. [19] devise heuristics that

leverage patterns in these interactions. They identify *round-trip transactions* where a unique amount enters the pool and leaves it soon after. We denote by  $w_{tz}$  a transaction from  $t$  to  $z$ , and  $w_{zt}$  from  $z$  to  $t$ . We represent the notion of time by a number  $t$  of transactions. We also define an equivalence relation  $V_a$  on  $\mathcal{S}$ , that denotes whether two sets of states  $\mathcal{S}_i, \mathcal{S}_j \in \mathcal{S}$  hold the same amount  $a$ , whenever  $\mathcal{S}_i V_a \mathcal{S}_j$ . This gives us the following inferred relation  $\tilde{L}$ :

$$\forall s_i \in \mathcal{S}_I \mathcal{E}_C w_{tz}, s_j \in w_{zt} \mathcal{E}_P \mathcal{S}_O, \exists! a : \quad (9)$$

$$\mathcal{S}_I V_a \mathcal{S}_O \wedge |\{w \mid s_i \prec w \prec s_j, w \in \mathcal{W}\}| \leq t \implies s_i \tilde{L} s_j$$

**Lemma 9.** *Any state  $s$  involved in a transaction  $T_w = \mathcal{S}_I w \mathcal{S}_O$  vulnerable to inferred relation  $\tilde{L}$  in (9) has an unlinkability score of  $u = \binom{|\mathcal{S}_I \cup \mathcal{S}_O \setminus [s]_{\tilde{L}}|}{2} / \binom{|\mathcal{S}_I \cup \mathcal{S}_O|}{2}$ .*

*Proof.* We have that  $[s]_{\tilde{L}}$  is the set  $\mathcal{S}_I \cup \mathcal{S}_O$ , for all  $s \in \mathcal{S}_I \cup \mathcal{S}_O$ . Therefore the unlinkability score of any state  $s$  is  $u = \binom{|\mathcal{S}_I \cup \mathcal{S}_O \setminus [s]_{\tilde{L}}|}{2} / \binom{|\mathcal{S}_I \cup \mathcal{S}_O|}{2}$ .

Furthermore, Kappos et al. [19] identify the retribution of mining pools towards their miners. If a transaction from the shielded pool to the transparent pool has many outputs states, one of which belongs to a known mining pool, all the states are from a mining pool. Let  $M$  be a refinement of the linking relation on the states  $\mathcal{S}$  that denotes whether two states  $s_i, s_j \in \mathcal{S}$  are miners, whenever  $s_i M s_j$ . We write  $[s]_M$  the equivalence class of  $\mathcal{S}$  by  $M$ , to which  $s$  belongs. Also, the threshold to the number of output is set to 100 by Kappos et al., here we generalize it to  $n$ . We express this heuristic as an inferred linking relation  $\tilde{L}$ :

$$\forall s_i, s_j \in w_{z-t} \mathcal{E}_P \mathcal{S}_O : s_i \in [s]_M \wedge |\mathcal{S}_O| > 100 \implies s_i \tilde{L} s_j \quad (10)$$

**Lemma 10.** *Any state  $s \notin [s']_M$  involved in a transaction  $T_w$  vulnerable to inferred relation  $\tilde{L}$  in (10) has an unlinkability score of 1, for all  $s' \in \mathcal{S}$ .*

*Proof.* If  $s$  is not in  $[s']_M$  for all  $s' \in \mathcal{S}$ , then  $[s']_M \cap [s]_{\tilde{L}} = \emptyset$ , since  $M$  is a refinement of  $L$ . Therefore, the score of  $s$  is  $\binom{|[s']_M|}{2} / \binom{|[s']_M|}{2} = 1$ .