RESEARCH ARTICLE

# Coin Transfer Unlinkability Under the Counterparty Adversary Model

Takeshi Miyamae,[*†] Kanta Matsuura[‡]

**Abstract.** Unlinkability is a crucial property of cryptocurrencies that protects users from deanonymization attacks. However, currently, even anonymous cryptocurrencies do not necessarily attain unlinkability under specific conditions. For example, Mimblewimble, which is considered to attain coin unlinkability using its transaction kernel offset technique, is vulnerable under the assumption that privacy adversaries can send their coins to or receive coins from the challengers. This paper first illustrates the privacy issue in Mimblewimble that could allow two colluded adversaries to merge a person's two independent chunks of personally identifiable information (PII) into a single PII. To analyze the privacy issue, we formulate *unlinkability* between two sets of objects and a privacy adversary model in cryptocurrencies called the *counterparty adversary model*. On these theoretical bases, we define an abstract model of blockchain-based cryptocurrency transaction protocols called the *coin transfer system*, and unlinkability over it called *coin transfer unlinkability (CT-unlinkability)*. Furthermore, we introduce zero-knowledgeness for the coin transfer systems to propose a method to easily prove the CT-unlinkability of cryptocurrency transaction protocols. Finally, we prove that Zerocash is CT-unlinkable by using our proving method to demonstrate its effectiveness.

## 1. Introduction

*1.1. PII Linkability Issue via Mimblewimble Protocol*—Unlinkability is a crucial property of cryptocurrencies that protects users from deanonymization attacks, as demonstrated by Bonneau *et al.*, Amarasinghe *et al.*, *etc.*;[1,2] however, although Silveira *et al.* derived the conclusion that Mimblewimble protocol attains transaction unlinkability, the actual risks of cryptocurrency's coin linkability are not necessarily understood.[3–5] In this section, we illustrate a privacy issue via Mimblewimble protocol concerning personally identifiable information (PII).[6]

Let "Alice" be an employee of Fujitsuna Co., and be assigned an employee number. Her official PII, *e.g.*, name, email address, phone number, address, ID photo, is recorded on the Fujitsuna Co.'s employee database. Because Fujitsuna Co. is a high-technology company, they pay the salaries of their employees using Beam,[7] one of the Mimblewimble-based cryptocurrencies.[3]

On the other hand, let Alice privately enjoy a content distribution service operated by Tsutayama Movie Contents (TMC) Co. She is assigned her user ID. Her private PII, *e.g.*, private

---

[*] 1ERsMcmcaRCri5jYPJ8Tqbr93V9UFXk3eF

[†] T. Miyamae (miyamae.takeshi@fujitsu.com) is a senior researcher at Fujitsu Limited, Japan and a Ph.D. student at The University of Tokyo, Japan.

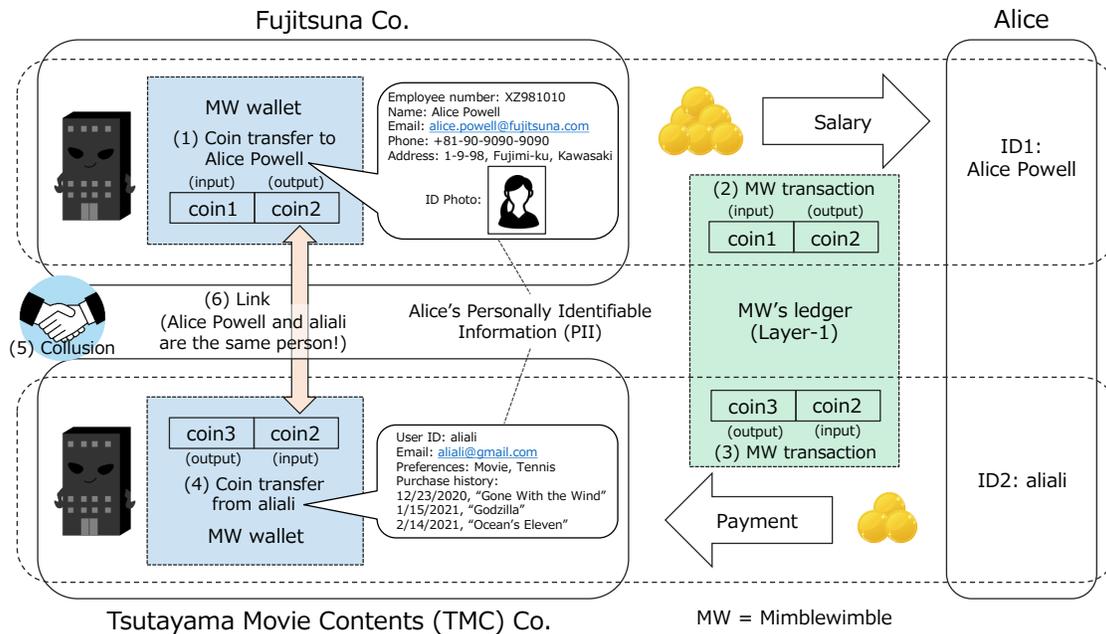[‡] K. Matsuura (kanta@iis.u-tokyo.ac.jp) is a professor at The University of Tokyo, Japan.

Fig. 1. PII Linkability Issue via Mimblewimble

email address, preferences including hobbies, and purchase history, is recorded on the TMC Co.'s customer database. Since TMC Co. is also a high-technology company, they support Beam as a payment means.

In most Mimblewimble-based cryptocurrencies, both the sender and the recipient of a transaction know all identifiers of input coins and output coins. As shown in Figure 1, Fujitsuna Co. and Alice know the identifiers of coin1 and coin2, while Alice and TMC Co. know the identifiers of coin2 and coin3. Therefore, if Fujitsuna Co. and TMC Co. colluded, they could associate Alice's official PII in Fujitsuna Co. with her private PII in TMC Co. via the identifier of coin2. That is, Fujitsuna Co. could grasp Alice's private information, *e.g.*, her hobbies and purchase history of movies, which is one of the privacy risks caused by Mimblewimble's insufficient coin unlinkability. However, this kind of issue has never been thoroughly discussed, and no solutions have ever been proposed so far. This is mainly because the sender and the recipient in a coin transfer are assumed to trust each other in most existing privacy adversary models.

In this study, we introduce the *counterparty adversary model* and *coin transfer unlinkability (CT-unlinkability)* to analyze the PII unlinkability via any kind of cryptocurrency transaction protocols.

*1.2.    Related Work*—Pfitzmann *et al.* systematically defined the unlinkability of two or more items of interest (IOIs) for general communication protocols.[8] Some cryptocurrency researchers, *e.g.*, Amarasinghe *et al.*, noticed that cryptocurrency transactions can also be handled in the same way as the messages in the communication protocols, and they attempted to apply the Pfitzmann's unlinkability to their analysis of cryptocurrency's privacy properties.[2] However, the application of the Pfitzmann's unlinkability does not seem reasonable because it is defined as a relation among the same kind of objects in a set (*e.g.*, message senders), while we would like to handle the relation between general objects (*e.g.*, a relation between a sender and a transaction).

Backes *et al.* proposed AnoA framework for defining, analyzing, and quantifying anonymity

properties, including *sender unlinkability*, for anonymous communication (AC) protocols, using computational differential privacy.[9] They modeled the AC protocols, generalized the notion of computational differential privacy (CDP) to apply it to the AC protocols, and analyzed the privacy of the AC protocols. However, since their theory assumes simple messaging that only includes a sender, a recipient, and auxiliary information (that corresponds to layer-0 in blockchain protocol layers), it is difficult for us to apply it to cryptocurrency transaction protocols (that correspond to layer-1 and layer-2).

In the research field of cryptocurrency, Androulaki *et al.* defined *activity unlinkability* and *user profile indistinguishability* to analyze the privacy of cryptocurrencies including Bitcoin.[10, 11] *Activity unlinkability* assesses the difficulty of identifying whether two different addresses or transactions belong to the same user. *User profile indistinguishability* assesses the adversary's ability to group the addresses and transactions of the same user. However, since their definitions only focus on the unlinkability (or indistinguishability) between addresses or transactions, they cannot handle unlinkability between general objects including coins. Furthermore, since it assumes that the challengers never share any secret information with the adversaries, it cannot be applied to the privacy issue under the *counterparty adversary model* defined in our work.

Ben-Sasson *et al.* defined *ledger indistinguishability* in their Zerocash paper to assess the unlinkability between the transactions on the ledger of their decentralized anonymous payment (DAP) scheme and the honest parties who participate in the DAP scheme.[12] The adversary can adaptively induce the honest parties to perform DAP operations of his choice. However, since they assume that the counterparties of each challenger are not included in their adversaries and that the private information of each coin is always hidden on the challenger's side, their privacy adversary model is not realistic. Furthermore, since the challenger's APIs are largely specific to the DAP scheme, their *ledger indistinguishability* is only applicable to Zerocash or its close relatives.

Silveira *et al.* defined *transaction unlinkability* in their formulation of Mimblewimble.[3] They claim using the definition that since transaction kernel offsets are added to generate a single block kernel offset, Mimblewimble attains *transaction unlinkability*. However, the scheme is specific to the confidential transactions and the claim is only applicable to Mimblewimble.[13]

In the research field of cryptography, a large number of zero-knowledge proof system protocols have been proposed based on the simulation paradigm.[14–16] These protocols have made a significant contribution to enhancing privacy in various fields including anonymous cryptocurrencies. However, the application of the simulation paradigm in any other field than zero-knowledge proof systems has been unusual so far.

*1.3. Contents*—The remainder of this paper is organized as follows. In Section 2, we formulate *unlinkability* between two sets of objects. In Section 3, we define a privacy adversary model in cryptocurrencies called the *counterparty adversary model*. In Section 4, we define an abstract model of blockchain-based cryptocurrency transaction protocols called the *coin transfer system*, and unlinkability over it called *coin transfer unlinkability (CT-unlinkability)* under the counterparty adversary model. In Section 5, we introduce zero-knowledgeness for the coin transfer systems to propose a method to easily prove the CT-unlinkability of cryptocurrency transaction protocols. In Section 6, we prove that Zerocash is CT-unlinkable by using our proving method to demonstrate its effectiveness. Section 7 concludes the paper.
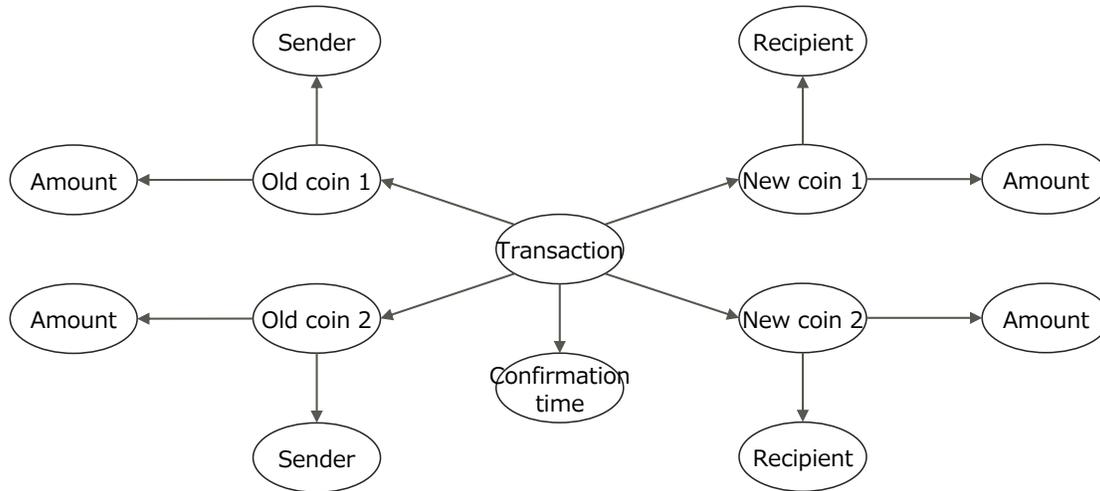
**19**

Fig. 2. Cryptocurrency Transaction Information Model

## 2.    Formulation of Unlinkability

This section formulates *unlinkability* between two sets of objects.

*2.1.    Cryptocurrency Transaction Information Model*—First, we introduce an object graph that represents a typical transaction information of UTXO-style cryptocurrency in Figure 2. Everyone instantly notices that the transaction sender and the transaction recipient are linkable if a cryptocurrency is designed to be perfectly transparent as Bitcoin is.  In contrast, most anonymous cryptocurrencies try to hide the link between the sender and the recipient.  If it is difficult to show the association between the sender and the recipient, we informally call it *unlinkable*.

*2.2.    Object*—We define a notion called *object* as the most fundamental element of the cryptocurrency transaction information model, over which we formally define several privacy properties in cryptocurrencies.  An object is defined as information because blockchain-based cryptocurrencies are implemented as software.

**Definition 1** (Object). *An 'object' is a primitive piece of information that indicates someone or something in a cryptocurrency.*

For example, a transaction, a sender, a recipient, an old coin, a new coin, coin amount, and transaction confirmation time are the objects in a typical UTXO-style cryptocurrency.

*2.3.    Attribute*—We define a notion called *attribute* to introduce a fundamental relation between objects.

**Definition 2** (Attribute). *An 'attribute' is an object that indicates a property of an original object. We also call the attribute a 'child' of the original object and call the original object a 'parent' of the attribute.*

For example, an old coin is an attribute of a transaction, and a sender is an attribute of an old coin in a typical UTXO-style cryptocurrency.
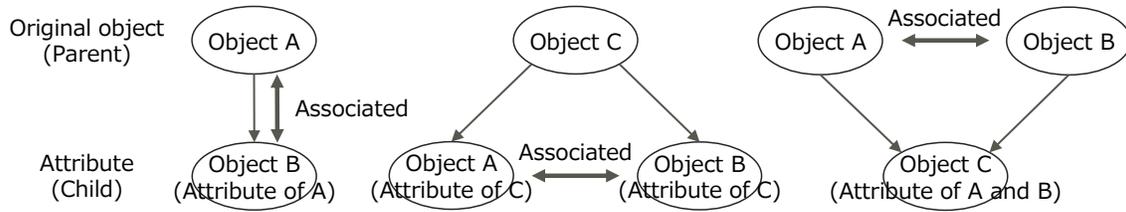
Fig. 3. Object, Attribute, and Association

*2.4. Association*—We define another notion of relation called *association* that does not contain any hierarchical relationship between objects.

**Definition 3** (Association). *Let $\boldsymbol{O}$ be a set of objects. If $A \in \boldsymbol{O}$ and $B \in \boldsymbol{O}$ satisfy either of the following conditions from the perspective of a participant (e.g., an adversary), we call the relation between A and B 'A and B are associated.' (We define a boolean function $Assoc : \boldsymbol{O}^2 \to \{0,1\}$ as $Assoc(A,B) = 1$ if and only if A and B are associated.)*

(1) *A is an attribute of B, or B is an attribute of A.*
(2) *Both A and B are attributes of an identical parent object.*
(3) *Both A and B have an identical child attribute object.*
(4) *There exists $C \in \boldsymbol{O}$ wherein $Assoc(A,C) = 1$ and $Assoc(B,C) = 1$.*

Note that the relation defined in Definition 3 is also expressed as '*A is associated with B*' or '*B is associated with A*.' If *A* and *B* are *not* associated, then we call the relation '*A and B are unassociated*.'

For example, a sender is associated with a sent coin because the sender is an attribute of the sent coin. A coin sent by a sender to Alice and a coin sent by the identical sender to Bob are associated because both the coin sent to Alice and the coin sent to Bob have an identical sender.

We illustrate these notions for the cryptocurrency transaction information model (object, attribute and association) in Figure 3.

*2.5. Unlinkability*—We define *unlinkability* in this section, which is the essential privacy property in cryptocurrencies.

**Definition 4** (Unlinkability). *Let $\boldsymbol{O}$ be a set of objects. Given a set of objects of an identical type $\boldsymbol{A} = \{A_i | i = 0,...,n-1\} \subset \boldsymbol{O}$ and another set of objects of another identical type $\boldsymbol{B} = \{B_j | j = 0,...,m-1\} \subset \boldsymbol{O}$, where $\boldsymbol{A} \cap \boldsymbol{B} = \varnothing$ .* If the adversary $\mathscr{A}$'s advantage of the indistinguishability game $\boldsymbol{G}_{assoc}(n_{sec})(\boldsymbol{A},\boldsymbol{B})$ where $n_{sec}$ is a security parameter (e.g. key length), $adv(\boldsymbol{G}_{assoc}(n_{sec})(\boldsymbol{A},\boldsymbol{B})) = max(|Pr[b = b']-1/2|)$, is negligible, then we call the relation between $\boldsymbol{A}$ and $\boldsymbol{B}$ '$\boldsymbol{A}$ and $\boldsymbol{B}$ are unlinkable.'*

*[Indistinguishability game $\boldsymbol{G}_{assoc}(n_{sec})(\boldsymbol{A},\boldsymbol{B})$]*

*Challenger $\mathscr{C}$: randomly selects a pair of objects $p_0 = (A_e \in \boldsymbol{A}, B_q \in \boldsymbol{B})$ where $Assoc(A_e,B_q) = 1$, randomly selects another pair of objects $p_1 = (A_f \in \boldsymbol{A}, B_s \in \boldsymbol{B})$ where $Assoc(A_f,B_s) = 0$, randomly selects $b \leftarrow \{0,1\}$, and sends $P_b$ to $\mathscr{A}$ according to the value of b where $P_0 = (p_0,p_1)$ and $P_1 = (p_1,p_0)$.*

*Adversary $\mathscr{A}$: guesses $b' \leftarrow \{0,1\}$ where $P_{b'}$ was sent by $\mathscr{C}$.*

---

\* We assume that at least one pair of objects are associated and at least another pair of objects are unassociated.
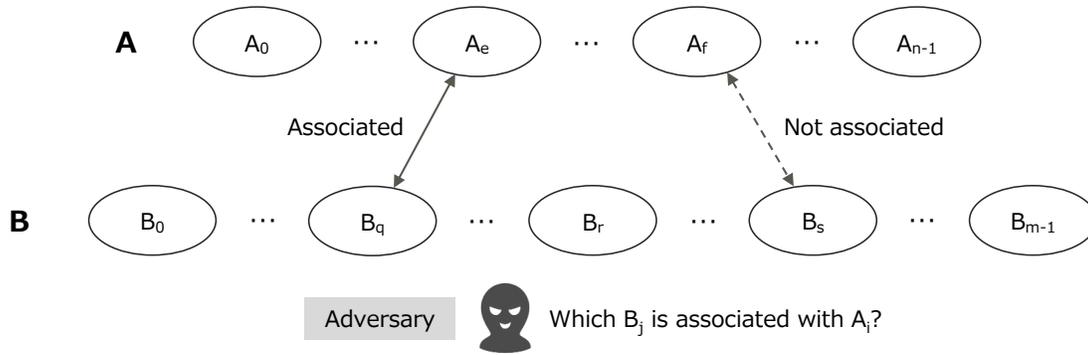
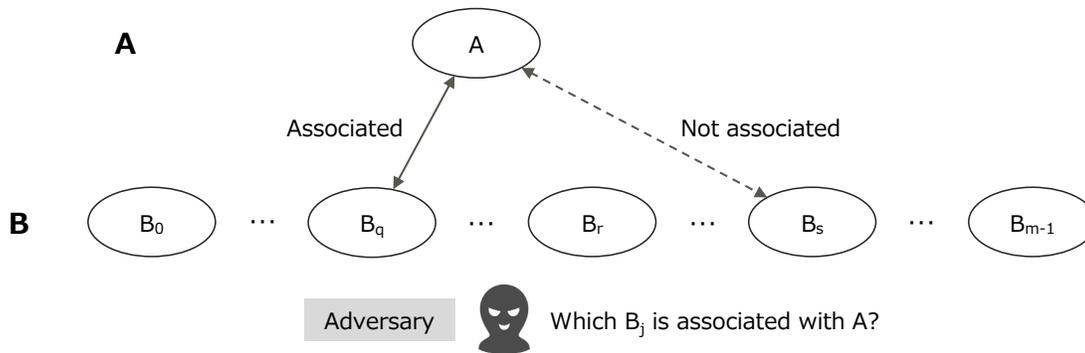Fig. 4. Indistinguishability Game $G_{assoc}(A, B)$ for Unlinkability



Fig. 5. Indistinguishability Game $G_{assoc}(A = \{A\}, B)$ for Object Unlinkability

Note that the relation defined in Definition 4 is also expressed as '$A$ *is unlinkable with* $B$' or '$B$ *is unlinkable with* $A$.' If $A$ and $B$ are *not* unlinkable, then we call the relation '$A$ *and* $B$ *are linkable.*'

Figure 4 depicts the indistinguishability game $G_{assoc}(A, B)$ for unlinkability.

For example, $B = \{cm_i \mid$ all the coin commitments on the ledger$\}$ is unlinkable with $A = \{addr_i \mid$ all the possible recipient addresses$\}$ in Zerocash because adversaries cannot guess who is the recipient of each coin commitment.[12]

On the other hand, $B = \{cm_i \mid$ all the coin descriptors (UTXO transaction outputs) on the ledger$\}$ is linkable with $A = \{addr_i \mid$ all the possible recipient addresses$\}$ in Bitcoin because every relation between a coin and its recipient is disclosed in the coin descriptor and anyone can see it.[11]

**Definition 5** (Object Unlinkability). *If $A$ contains only a single object A ($A = \{A\}$), then Definition 4 defines 'object unlinkability' between $A = \{A\}$ and $B$ using the indistinguishability game $G_{assoc}(A = \{A\}, B)$.*

Figure 5 depicts the indistinguishability game $G_{assoc}(A = \{A\}, B)$ for object unlinkability.

For example, $A$ (= a recipient address) is unlinkable with $B = \{cm_i \mid$ all the coin commitments on the ledger$\}$ in Zerocash. Conversely, $A$ (= a coin commitment) is unlinkable with $B = \{addr_i \mid$ all the possible recipient addresses$\}$ in Zerocash.

**Theorem 2.1** (Transitivity of Unlinkability). *Unlinkability is a transitive relation.*

[ Proof. ] Let each of $\boldsymbol{A}$ and $\boldsymbol{B}$ be a set of objects of an identical type and $\boldsymbol{A}$ and $\boldsymbol{B}$ are unlinkable. The probability that an adversary wins the indistinguishability game $\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{B})$ is

$$
\begin{aligned}
Pr(\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{B}))[b=b'] &= Pr[(b=0 \wedge b'=0) \vee (b=1 \wedge b'=1)] \\
&= L(A_e,B_q) = \frac{l(A_e,B_q)}{l(A_e,B_q)+l(A_f,B_s)} \\
&< \frac{1}{2} + adv(\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{B}))
\end{aligned}
$$

that is,

$$
max(l(A_i,B_j)) < \frac{1+2adv(\boldsymbol{A},\boldsymbol{B})}{1-2adv(\boldsymbol{A},\boldsymbol{B})} min(l(A_i,B_j))
$$

where $l(A_i,B_j)$ is the likelihood that $B_j$ is associated with $A_i$ from the adversary's perspective, and $L(A_i,B_j)$ is the conditional likelihood that $B_j$ is associated with $A_i$ from the adversary's perspective after they are sent $Pb$ in the indistinguishability game $\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{B})$.

Likewise, if $\boldsymbol{C}$ is another set of objects of another identical type and $\boldsymbol{B}$ and $\boldsymbol{C}$ are unlinkable, the probability that an adversary wins the indistinguishability game $\boldsymbol{G}_{assoc}(\boldsymbol{B},\boldsymbol{C})$ is

$$
\begin{aligned}
Pr(\boldsymbol{G}_{assoc}(\boldsymbol{B},\boldsymbol{C}))[b=b'] &= Pr[(b=0 \wedge b'=0) \vee (b=1 \wedge b'=1)] \\
&= L(B_q,C_v) = \frac{l(B_q,C_v)}{l(B_q,C_v)+l(B_s,C_w)} \\
&< \frac{1}{2} + adv(\boldsymbol{G}_{assoc}(\boldsymbol{B},\boldsymbol{C}))
\end{aligned}
$$

that is,

$$
max(l(B_j,C_k)) < \frac{1+2adv(\boldsymbol{B},\boldsymbol{C})}{1-2adv(\boldsymbol{B},\boldsymbol{C})} min(l(B_j,C_k))
$$

We now consider the association between $A_i \in \boldsymbol{A}$ and $C_k \in \boldsymbol{C}$, deduced from the transitivity of association, and the indistinguishability game $\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{C})$ (see Figure 6). The probability that an adversary wins the indistinguishability game $\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{C})$ is

$$
\begin{aligned}
Pr(\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{C}))[b=b'] &= Pr[(b=0 \wedge b'=0) \vee (b=1 \wedge b'=1)] \\
&= L(A_e,C_v) = \frac{l(A_e,C_v)}{l(A_e,C_v)+l(A_f,C_w)} \\
&= \frac{\sum_{B_j} l(A_e,B_j)l(B_j,C_v)}{\sum_{B_j} l(A_e,B_j)l(B_j,C_v) + \sum_{B_j} l(A_f,B_j)l(B_j,C_w)}
\end{aligned}
$$

$$
\begin{aligned}
&< \frac{|B|\frac{1+2adv(\boldsymbol{A},\boldsymbol{B})}{1-2adv(\boldsymbol{A},\boldsymbol{B})}min(l(A_i,B_j))\frac{1+2adv(\boldsymbol{B},\boldsymbol{C})}{1-2adv(\boldsymbol{B},\boldsymbol{C})}min(l(B_j,C_k))}{|B|\frac{1+2adv(\boldsymbol{A},\boldsymbol{B})}{1-2adv(\boldsymbol{A},\boldsymbol{B})}min(l(A_i,B_j))\frac{1+2adv(\boldsymbol{B},\boldsymbol{C})}{1-2adv(\boldsymbol{B},\boldsymbol{C})}min(l(B_j,C_k)) + |B|min(l(A_i,B_j))min(l(B_j,C_k))} \\
&= \frac{\{1+2adv(\boldsymbol{A},\boldsymbol{B})\}\{1+2adv(\boldsymbol{B},\boldsymbol{C})\}}{\{1+2adv(\boldsymbol{A},\boldsymbol{B})\}\{1+2adv(\boldsymbol{B},\boldsymbol{C})\} + \{1-2adv(\boldsymbol{A},\boldsymbol{B})\}\{1-2adv(\boldsymbol{B},\boldsymbol{C})\}} \\
&= \frac{1+2adv(\boldsymbol{A},\boldsymbol{B}) + 2adv(\boldsymbol{B},\boldsymbol{C}) + 4adv(\boldsymbol{A},\boldsymbol{B})adv(\boldsymbol{B},\boldsymbol{C})}{2+8adv(\boldsymbol{A},\boldsymbol{B})adv(\boldsymbol{B},\boldsymbol{C})} \\
&\to \frac{1}{2} + \{adv(\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{B})) + adv(\boldsymbol{G}_{assoc}(\boldsymbol{B},\boldsymbol{C}))\} = \frac{1}{2} + adv(\boldsymbol{G}_{assoc}(\boldsymbol{A},\boldsymbol{C}))
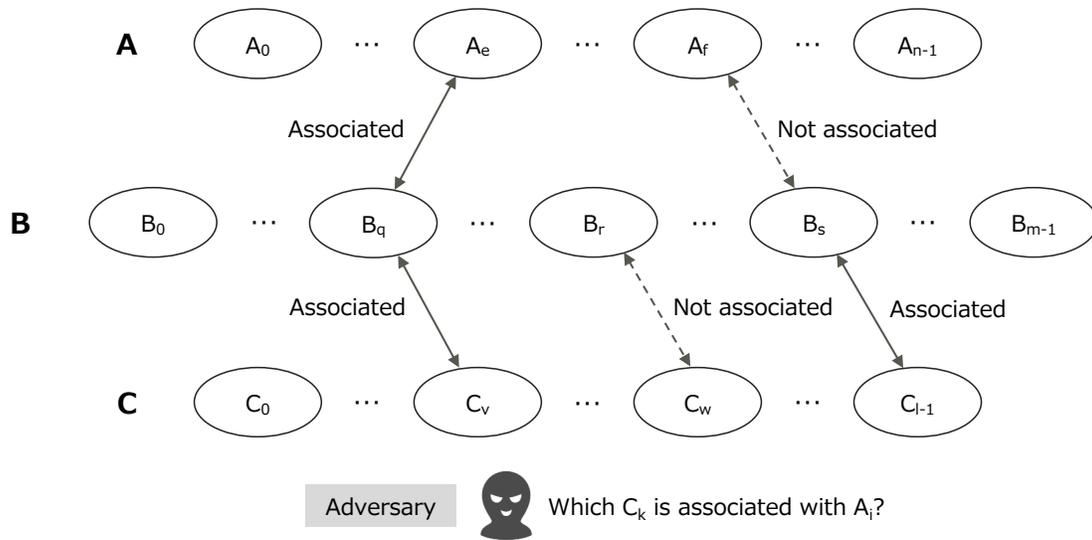\end{aligned}
$$

Fig. 6. Indistinguishability Game $G_{assoc}(A, C)$ to Prove the Transitivity of Unlinkability

Hence, $A$ and $C$ are unlinkable (unlinkability is a transitive relation). Q.E.D.

From this result, we only have to focus on the unlinkability in each information layer to achieve the entire unlinkability over the whole cryptocurrency transaction information model.

## 3. Privacy Adversary Models

In this section, we define a privacy adversary model in cryptocurrencies called the *counterparty adversary model*, where we assume that the counterparties of each challenger in coin transfer transactions can become their privacy adversaries.

*3.1. Current Privacy Adversary Model*—Most of the past cryptocurrency studies have assumed that the counterparties of any challengers in coin transfer transactions never become their privacy adversaries.

The following two simple rules are assumed in the current privacy adversary model (see Figure 7).

(1) The adversaries may not send their coins to the challengers.
(2) The challengers may not send their coins to the adversaries.

As a result, the adversaries can only observe the coin transfer transactions between the challengers as privacy attacks. They cannot share any secret transaction information with the challengers.

*3.2. Counterparty Adversary Model*—On the other hand, we propose another privacy adversary model in cryptocurrencies, which does not assume either rule defined in Section 3.1 (see Figure 8). We call it the *counterparty adversary model*. In this privacy adversary model, since the adversaries can send their coins to or receive coins from the challengers, they can observe the challengers' privacy using their shared secrets, in addition to the coin transfer observation.

The counterparty adversary model is highly reasonable because the real world is not split cleanly into two domains as shown in the current privacy adversary model. However, the counterparty adversary model has never been thoroughly discussed so far. Therefore, we must
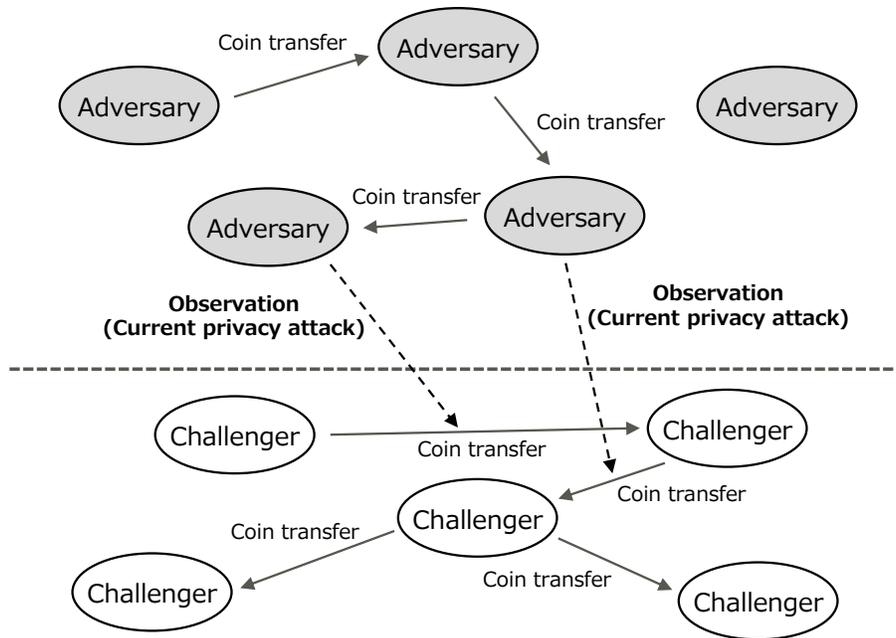
Fig. 7. Current Privacy Adversary Model

estimate the privacy threat from the counterparties during coin transfers to use cryptocurrencies as everyday payment means.

*3.3. PII Unlinkability*—For example, we define a privacy issue that only occurs under the counterparty adversary model, called *PII unlinkability* (Figure 9). We consider a case that an adversary $\mathscr{A}_1$ sends a coin to a challenger $\mathscr{C}_1$, and $\mathscr{C}_1$ sends a coin using another identity $\mathscr{C}_2$ to another adversary $\mathscr{A}_2$. Even if $\mathscr{A}_1$ and $\mathscr{A}_2$ colluded with each other and they shared the $\mathscr{C}_1$'s PII and the $\mathscr{C}_2$'s PII with each other, a perfect privacy-enhancing coin transfer would keep the shared these PIIs unlinkable.[†] Conversely, if $\mathscr{A}_1$ and $\mathscr{A}_2$ colluded with each other and the coins that $\mathscr{A}_1$ sent to $\mathscr{C}_1$ and the coins that $\mathscr{C}_2$ sent to $\mathscr{A}_2$ were linkable, the adversaries could conclude that both identities $\mathscr{C}_1$ and $\mathscr{C}_2$ belong to an identical person. They could merge the $\mathscr{C}_1$'s PII and the $\mathscr{C}_2$'s PII into a single PII.

## 4. Coin Transfer System

In this section, we define an abstract model of blockchain-based cryptocurrency transaction protocols called the *coin transfer system*, and unlinkability over it called *coin transfer unlinkability (CT-unlinkability)* under the counterparty adversary model.

**Definition 6** (Coin Transfer System). *Let $\boldsymbol{M}$ be a set of probabilistic polynomial-time Turing machines. $\boldsymbol{M}$ is a 'coin transfer system' if every machine $M_{from} \in \boldsymbol{M}$ can transfer their ownership of the coin to any other machine $M_{to} \in \boldsymbol{M}$ using a public information channel and a private information channel between them. We call each $M \in \boldsymbol{M}$ a 'coin transfer machine.'*

Note that a *probabilistic polynomial-time algorithm* means a probabilistic algorithm that always (*i.e.*, independently of the outcome of its internal coin tosses) halts after a polynomial (in

---

[†] We assume that the $\mathscr{C}_1$'s PII and the $\mathscr{C}_2$'s PII have no common part.
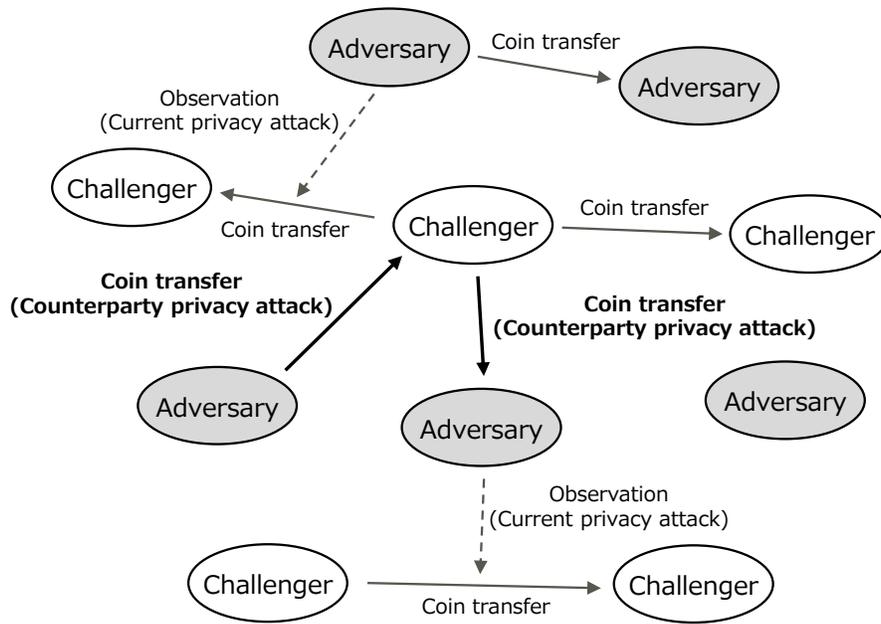
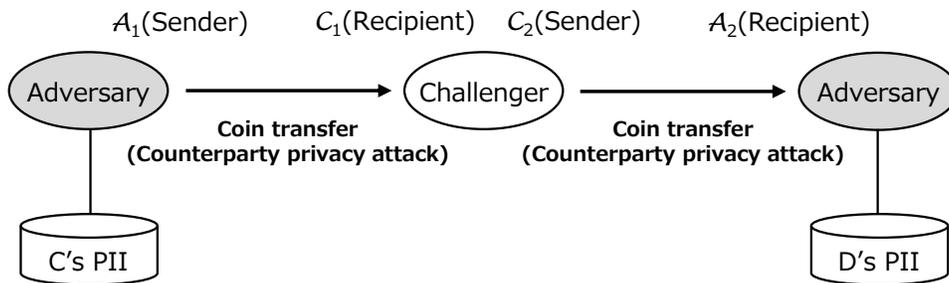Fig. 8. Counterparty Adversary Model
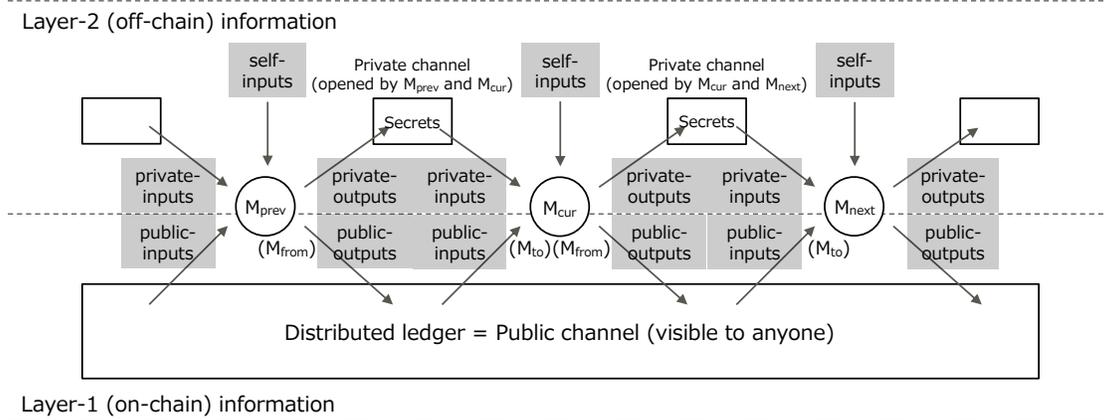


Fig. 9. PII Unlinkability

**26**

Fig. 10. Coin Transfer System

the length of the input) number of steps. The adversary's verdict of whether to accept or reject the *unlinkabilitiy* is probabilistic, and a bounded error probability is allowed.

For example, a coin transfer system for blockchain-based cryptocurrencies is illustrated in Figure 10. Every blockchain-based cryptocurrency, including Bitcoin, uses a distributed ledger as a public information channel. Furthermore, most anonymous cryptocurrencies provide their protocols, including private channels between each participant.

Focusing on a machine $M_{cur}$, which is sent coins from $M_{prev}$ and sends the sent coins to $M_{next}$, its inputs are classified into three types, public-inputs (from the ledger), private-inputs (from $M_{prev}$), and self-inputs (*e.g.*, random value generated by $M_{cur}$). On the other hand, its outputs are classified into two types, public-outputs (to the ledger) and private-outputs (to $M_{next}$). This transaction process is denoted as a function named $M_{cur}$, such as (public-outputs, private-outputs to $M_{next}$) = $M_{cur}$(public-inputs, private-inputs from $M_{prev}$, self-inputs).

Next, we introduce unlinkability in coin transfer systems under the counterparty adversary model.

**Definition 7** (Coin Transfer Unlinkability). *Given a coin transfer system defined in Definition 6 and a pair of colluding adversaries $M_{prev}$ and $M_{next}$. If all the coins sent by $M_{prev}$ to the addresses other than $M_{prev}$ and $M_{next}$, denoted as $\boldsymbol{C}_{out}$, and all the coins sent by the addresses other than $M_{prev}$ and $M_{next}$ to $M_{next}$, denoted as $\boldsymbol{C}_{in}$, are unlinkable (regarding the association such that one coin is spent by another), we call the relation between $\boldsymbol{C}_{out}$ and $\boldsymbol{C}_{in}$ 'coin transfer unlinkable (CT-unlinkable).'*

The CT-unlinkability defined in Definition 7 is essential under the counterparty adversary model because it ensures PII unlinkability shown in Figure 9.

**Theorem 4.1** (PII unlinkability under CT-unlinkability). *CT-unlinkability ensures PII unlinkability.*

[ Proof. ] Assuming that $\boldsymbol{C}_{out}$ and $\boldsymbol{C}_{in}$ are CT-unlinkable and that the $\mathscr{C}_1$'s PII and the $\mathscr{C}_2$'s PII are linkable by any coin transfers, at least one coin sent by adversary $\mathscr{A}_1$ to $\mathscr{C}_1$ and another coin sent by $\mathscr{C}_2$ to adversary $\mathscr{A}_2$ seems to be associated, which means that $\boldsymbol{C}_{out}$ and $\boldsymbol{C}_{in}$ are CT-linkable. This situation is a contradiction, so the assumption must be false. Therefore, if $\boldsymbol{C}_{out}$

**27**

and $\boldsymbol{C}_{in}$ are CT-unlinkable, the $\mathscr{C}_1$'s PII and the $\mathscr{C}_2$'s PII are also unlinkable by any coin transfers. Q.E.D.

## 5.  Zero-Knowledge Coin Transfer System

This section introduces zero-knowledgeness for the coin transfer systems to propose a method to easily prove the CT-unlinkability of cryptocurrency transaction protocols.

It is relatively easy to prove that a certain cryptocurrency is CT-linkable because in such cases we only have to show at least one counterexample in which the adversary wins. In contrast, it is much more difficult to prove CT-unlinkable because we have to show that there are no examples. To show the CT-unlinkability easily, we introduce *zero-knowledgeness* using the *simulation paradigm* for the coin transfer systems similarly for zero-knowledge proof systems.[17]

**Definition 8** (Computational Zero-Knowledge Coin Transfer System). *Let $\boldsymbol{M}$ be a coin transfer system, $L_{pub}$ be the set of all the possible public-inputs of $\boldsymbol{M}$, and $L_{prv}$ be the set of all the possible private-inputs of $\boldsymbol{M}$. We say that $\boldsymbol{M}$ is 'computational zero-knowledge' (or just 'zero-knowledge') if for every probabilistic polynomial-time coin transfer machines $M_{prev}$, $M_{cur}$, $M_{next} \in \boldsymbol{M}$ there exists a probabilistic polynomial-time algorithm $M^*$ that does not use private-inputs from $M_{prev}$, such that the following two ensembles are computationally indistinguishable:*
*   $\{(public\text{-}outputs, private\text{-}outputs\ to\ M_{next}) = M_{cur}(public\text{-}inputs, private\text{-}inputs\ from\ M_{prev}, self\text{-}inputs)\}_{public\text{-}inputs \in L_{pub}, private\text{-}inputs \in L_{prv}}$ *(i.e., the output of the coin transfer machine $M_{cur}$ after receiving coins from the coin transfer machine $M_{prev}$ on public-inputs)*
*   $\{(public\text{-}outputs, private\text{-}outputs\ to\ M_{next}) = M^*(public\text{-}inputs, private\text{-}inputs\ from\ M_{prev}, self\text{-}inputs)\}_{public\text{-}inputs \in L_{pub}, private\text{-}inputs \in L_{prv}}$ *(i.e., the output of the algorithm $M^*$ on public-inputs)*
*Algorithm $M^*$ is called a 'computational simulator' (or just 'simulator') for the transformation $M_{cur}$ of the private-inputs from $M_{prev}$.*

The simulation paradigm works as a proving method to show that a participant has never been informed of any knowledge from the protocol. In the case of Definition 8, the existence of the simulator $M^*$ shows that no new association knowledge between cryptocurrency objects is leaked from $M_{cur}$ to $M_{next}$.

We also define an additional property called *universality* for zero-knowledge coin transfer systems.

**Definition 9** (Universal Zero-Knowledge Coin Transfer System). *Let $\boldsymbol{M}$ be a (computational) zero-knowledge coin transfer system. Suppose a simulator $M_u^*$ is only dependent on its self-inputs, and thus can be used for all $M_{cur} \in \boldsymbol{M}$. In that case, we call $M_u^*$ a universal simulator, and we call $\boldsymbol{M}$ a universal zero-knowledge coin transfer system.*

Here, we obtain a theorem that easily proves that a certain cryptocurrency is CT-unlinkable.

**Theorem 5.1** (CT-Unlinkability of Universal Zero-knowledge Coin Transfer Systems). *Universal zero-knowledge coin transfer systems are CT-unlinkable.*

[ Proof. ] Let $\boldsymbol{M}$ be a universal zero-knowledge coin transfer system, $M_{prev}, M_{next} \in \boldsymbol{M}$ be the adversaries of the CT-unlinkability indistinguishability game who collude with each other, and thus $\{M_i | M_i \in \boldsymbol{M}, M_i \neq M_{prev}, M_i \neq M_{next}\}$ be the challengers of this game. Let a coin $c_{assoc} \in \boldsymbol{C}_{in}$ that is associated with a coin sent by $M_{prev}$ be sent by $M_{assoc}$, and a coin $c_{nassoc} \in \boldsymbol{C}_{in}$ that is not

associated with any coin sent by $M_{prev}$ be sent by $M_{nassoc}$. Since **M** is universally zero-knowledge, $M_{assoc}$ and $M_{nassoc}$ can use the same simulator $M_u^*$ that is only dependent on its self-inputs. Therefore, the outputs of $M_{assoc}$ and the outputs of $M_u^*$ are indistinguishable, while the outputs of $M_{nassoc}$ and the outputs of $M_u^*$ are also indistinguishable. Therefore, the outputs of $M_{assoc}$ and the outputs of $M_{nassoc}$ are indistinguishable, meaning that the adversaries cannot win the CT-unlinkability indistinguishability game. Thus **M** is CT-unlinkable. Q.E.D.

## 6. CT-unlinkability of Existing Cryptocurrency Transaction Protocols

In this section, we first prove that Mimblewimble is CT-linkable. Next, we prove that Zerocash is CT-unlinkable by using our proving method to demonstrate its effectiveness.

*6.1. Mimblewimble*—Mimblewimble is one of the major anonymous cryptocurrency transaction protocols that has already been implemented in Grin and Beam, and later formulated by Silveira *et al.*[3–5,7,18] Owing to the transaction kernel offset, Mimblewimble can shuffle the transaction outputs and attain coin unlinkability, as shown in ValueShuffle and the Silveira *et al.*'s formulation (ValueShuffle also proposes a decentralized protocol using secure multi-party computation).[19] Therefore, it is difficult for the adversaries under the current privacy adversary model to find any association between input coins (*e.g.*, coin1 in Figure 11) and output coins (coin2).

However, under the counterparty adversary model, a sender adversary $M_{prev}$ can know the association between the input coins (coin1) sent by the adversary $M_{prev}$ and the output coins (coin2) received by a challenger $M_{cur}$. Similarly, a receiver adversary $M_{next}$ can know the association between the input coins (coin2) sent by a challenger $M_{cur}$ and the output coins (coin3) received by the adversary $M_{next}$. Therefore, in the CT-unlinkability indistinguishability game adversaries can easily distinguish whether each coin sent by the sender adversary was spent for another coin received by the receiver adversary or not, which means that the adversaries' advantage in the indistinguishability game is not negligible. Thus, we conclude that Mimblewimble is *CT-linkable*.
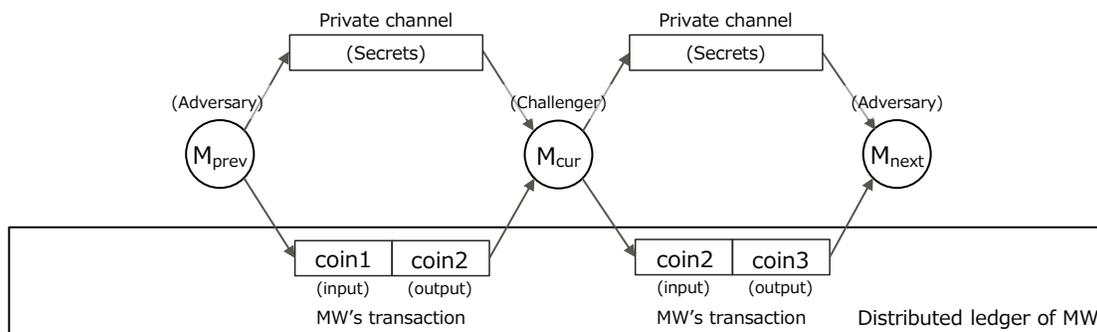


Fig. 11. Mimblewimble Coin Transfer System

*6.2. Zerocash*—Zerocash is one of the most sophisticated anonymous cryptocurrencies and has already been implemented in Zcash.[12,20] First, we briefly illustrate the anonymization scheme in Zerocash.

In Zerocash, all coins on the ledger are encrypted. Two seemingly unrelated identifiers, a coin commitment *cm* and a serial number *sn* (Figure 12), are used to identify each coin. A transaction
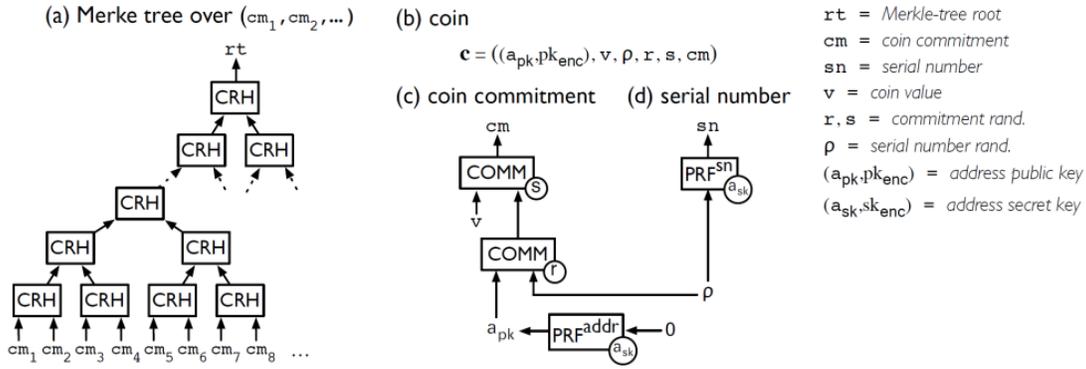
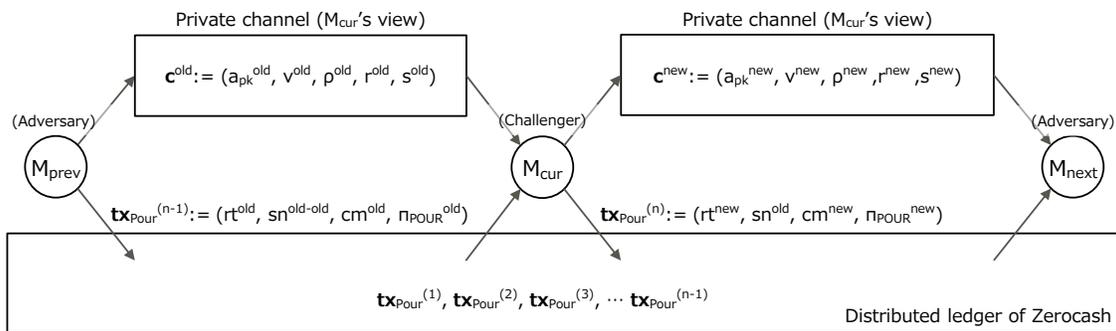Fig. 12. Zerocash Coin Commitment (cited from the Zerocash paper[12])



Fig. 13. Zerocash Coin Transfer System

sender uses $cm^{new}$ as an identifier of the coin (described as $c^{new}$ from the viewpoint of the sender). In contrast, a corresponding transaction recipient uses $sn^{old}$ as another identifier of the same coin (described as $c^{old}$ from the viewpoint of the recipient). Consequently, this makes the two identifiers of the same coin recorded in the shared ledger unlinkable even from the sender (only the recipient can link those two coin identifiers).

To convince the recipient that the sent coins are valid (*i.e.*, the shared secrets sent from the sender are well-formed), the sender includes a zero-knowledge proof $\pi_{POUR}^{new}$ in each transaction to prove that the identifiers and the secret parameters of sent coins are consistent, as shown in the followings (the suffixes are added from the viewpoint of the sender).

- The coin commitment $cm^{old}$ is properly calculated from the secret information of the coin $c^{old}$ sent to and held by the sender.
- The coin commitment $cm^{new}$ is properly calculated from the secret information of the coin $c^{new}$ produced by the sender and sent to the recipient.
- The public key $a_{pk}^{old}$ used as an address when the sender has received $cm^{old}$ is properly generated from the sender's secret key $a_{sk}^{old}$ (the proof of the proper recipient $a_{pk}^{old}$ of $c^{old}$).
- The serial number $sn^{old}$ is properly calculated from the secret information of the coin $c^{old}$ held by the sender and the sender's secret key $a_{sk}^{old}$.
- The secret information $cm^{old}$ is included in the CMList $rt^{new}$ recorded in the shared ledger

**30**

(the proof of the fact that $cm^{old}$ had already existed when the sender received the coin $\boldsymbol{c}^{old}$).

- The total amount of old coins ($v_1^{old} + v_2^{old} + ...$) included in the transaction input is equal to the total amount of new coins ($v_1^{new} + v_2^{new} + ...$) included in the transaction output.

Preventing double-spending of the coin $\boldsymbol{c}^{old}$ does not need to be proved using zero-knowledge proofs. It is sufficient to verify that $sn^{old}$ has not yet appeared in any transaction input recorded in the shared ledger.

Next, we prove that Zerocash is a universal zero-knowledge coin transfer system in Theorem 6.1.

**Theorem 6.1.** *Zerocash is a universal zero-knowledge coin transfer system.*

[ Proof. ] Let $\boldsymbol{M}_z$ be a coin transfer system of Zerocash. We show that, for every probabilistic polynomial-time participants $M_{prev}, M_{cur}, M_{next} \in \boldsymbol{M}_z$, we can construct a universal simulator $M_u^*$ for the Zerocash's POUR transaction, where the ensemble of $M_{cur}$'s outputs and the ensemble of $M_u^*$'s outputs are computationally indistinguishable, and the identical simulator $M_u^*$ that is only dependent on its self-inputs can be used for all $M_{cur} \in \boldsymbol{M}_z$.

Note that we do not construct the simulator regarding the public output parameters and implementation-dependent parameters (see Zerocash paper[12] for more details).

(1) Public-outputs of $M_u^*$

- CMList $rt^{sim} = CRH(rd^{sim})$; $rd^{sim} = rand() \in$ self-inputs of $M_u^*$ (Note that rand() is a pseudorandom number generator.)
  CMList $rt^{new}$ of $M_{cur}$'s public-outputs is a Merkle-tree root and thus a returned value of a collision-resistant hash function $CRH(x)$ wherein $x = cc_{root} = concat(CRH(cc_0), CRH(cc_1))$ ($cc_0$ and $cc_1$ are other concatenated returned values of $CRH(x)$). Since $CRH(x)$ is uniform, $rt^{new} = CRH(cc_{root})$ is computationally indistinguishable from $rt^{sim} = CRH(rd^{sim})$ that is only dependent on the self-inputs of $M_u^*$. Therefore, $M_u^*$ can universally simulate the public-output $rt^{new}$ of $M_{cur}$.
- Serial number $sn^{old\,sim} = PRF[a_{sk}{}^{old\,sim}](\rho^{old\,sim})$; $a_{sk}{}^{old\,sim}, \rho^{old\,sim} = rand() \in$ self-inputs of $M_u^*$
  Serial number $sn^{old}$ of $M_{cur}$'s public-outputs is a returned value of a pseudo-random function $PRF[a_{sk}{}^{old}](x)$ wherein $x = \rho^{old}$ (a random value generated by $M_{prev}$), which is computationally indistinguishable from $PRF[a_{sk}{}^{old}](x)$ wherein $x = \rho^{old\,sim} = rand()$ (a random value generated by $M_u^*$ as shown below). Furthermore, $PRF[a_{sk}{}^{old}](\rho^{old})$ is also computationally indistinguishable from $sn^{old\,sim} = PRF[a_{sk}{}^{old\,sim}](\rho^{old\,sim})$ where $M_u^*$ replaces $a_{sk}{}^{old}$ with $a_{sk}{}^{old\,sim}$ (a secret key of $M_u^*$). Since $sn^{old\,sim}$ is only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the public-output $sn^{old}$ of $M_{cur}$.
- Coin commitment $cm^{sim} = COMM_{s^{sim}}(v^{sim}||COMM_{r^{sim}}(a_{pk}{}^{sim}||\rho^{sim}))$; $a_{pk}{}^{sim} = a_{pk}{}^{new}$, $v^{sim} = v^{new}, \rho^{sim} = \rho^{new}, r^{sim} = r^{new}, s^{sim} = s^{new} \in$ self-inputs of $M_u^*$
  Coin commitment $cm^{new}$ of $M_{cur}$'s public-outputs is a returned value of a commitment function $COMM_x(z)$ wherein $x = s^{new}$ denotes the commitment trapdoor, and $z = v^{new}||COMM_{r^{new}}(a_{pk}{}^{new}||\rho^{new})$ denotes the committed value. $M_u^*$ can use the same coin commitment $cm^{sim} = cm^{new}$ as $M_{cur}$, because $a_{pk}{}^{sim} = a_{pk}{}^{new}, v^{sim} = v^{new}, \rho^{sim} = \rho^{new}$,

$r^{sim} = r^{new}$, $s^{sim} = s^{new}$, as shown below. Since $cm^{sim}$ is only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the private-output $cm^{new}$ of $M_{cur}$.

- ZKP $Sim[\pi_{POUR}{}^{sim}] = Sim[Prove(pk_{POUR}, \boldsymbol{x} = (rt^{old\,sim}, sn^{old\,sim}, cm^{sim}), \boldsymbol{a})]$; $rt^{old\,sim}$, $sn^{old\,sim}$, $cm^{sim} \in$ self-inputs of $M_u^*$

  ZKP $\pi_{POUR}{}^{new}$ of $M_{cur}$'s public-outputs is a returned value of a zk-SNARK prove function $Prove(pk_{POUR}, \boldsymbol{x}, \boldsymbol{a})$ wherein $\boldsymbol{x} = (rt^{old}, sn^{old}, cm^{new})$ denotes an instance of the NP language $\mathscr{L}$ of the statement $POUR$, and $\boldsymbol{a} = (a_{sk}{}^{old}, \boldsymbol{c}^{old}, \boldsymbol{c}^{new})$ denotes a witness. Assuming that $M_u^*$ send a coin $\boldsymbol{c}^{old\,sim} = (a_{pk}{}^{old\,sim}, v^{old\,sim}, \rho^{old\,sim}, r^{old\,sim}, s^{old\,sim}, cm^{old\,sim})$, where $a_{pk}{}^{old\,sim}$ is a public key of $M_u^*$, $v^{old\,sim}$ is larger than any assumed $v^{new}$, $\rho^{old\,sim} = rand()$, $r^{old\,sim} = rand()$, and $s^{old\,sim} = rand()$, to itself in the past (before $rt^{old\,sim}$). Under this assumption, $M_u^*$ can generate a proof $\pi_{POUR}{}^{sim} = Prove(pk_{POUR}, \boldsymbol{x} = (rt^{old\,sim}, sn^{old\,sim}, cm^{sim}), \boldsymbol{a})$ that meets completeness of zero-knowledge proof. $\pi_{POUR}{}^{new}$ and $\pi_{POUR}{}^{sim}$ are computationally indistinguishable because their inputs are mapped to a set of exponent values in the proof of zk-SNARK proposed by Parno *et al.*[14] Furthermore, $\pi_{POUR}{}^{sim}$ is also computationally indistinguishable from its zero-knowledge proof simulator $Sim[\pi_{POUR}{}^{sim}]$ that is not dependent on $\boldsymbol{a}$. Since $Sim[\pi_{POUR}{}^{sim}]$ is only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the private-output $\pi_{POUR}{}^{new}$ of $M_{cur}$.

(2) Private-outputs of $M_u^*$

- Address public key $a_{pk}{}^{sim} = a_{pk}{}^{new} \in$ self-inputs of $M_u^*$
  Address public key $a_{pk}{}^{new}$ of $M_{cur}$'s private-outputs is a number. $M_u^*$ can use the same address public key $a_{pk}{}^{sim} = a_{pk}{}^{new}$ as $M_{cur}$. Since $a_{pk}{}^{sim}$ is only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the private-output $a_{pk}{}^{new}$ of $M_{cur}$.

- Coin value $v^{sim} = v^{new} \in$ self-inputs of $M_u^*$
  Coin value $v^{new}$ of $M_{cur}$'s private-outputs is a number. $M_u^*$ can use the same coin value $v^{sim} = v^{new}$ as $M_{cur}$. Since $v^{sim}$ is only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the private-output $v^{new}$ of $M_{cur}$.

- Randomly sampled values $\rho^{sim} = \rho^{new} = rand()$, $r^{sim} = r^{new} = rand()$, $s^{sim} = s^{new} = rand() \in$ self-inputs of $M_u^*$
  Randomly sampled values $\rho^{new}$, $r^{new}$, $s^{new}$ of $M_{cur}$'s private-outputs are the random values generated by $M_{cur}$. $M_u^*$ can use the same random values $\rho^{sim} = \rho^{new}$, $r^{sim} = r^{new}$, $s^{sim} = s^{new}$ as $M_{cur}$. Since $\rho^{sim}, r^{sim}, s^{sim}$ are only dependent on the self-inputs of $M_u^*$, $M_u^*$ can universally simulate the private-outputs $\rho^{new}$, $r^{new}$ $s^{new}$ of $M_{cur}$.

Q.E.D.

Therefore, we finally conclude that Zerocash is *CT-unlinkable* from Theorem 5.1 and Theorem 6.1.

## 7. Conclusion

Unlinkability is a crucial property of cryptocurrencies that protects users from de-anonymization attacks. This paper first illustrated a privacy issue in Mimblewimble that could allow two colluded adversaries to merge a person's two independent chunks of personally identifiable information (PII) into a single PII.

To analyze the privacy issue, we formulated *unlinkability* between two sets of objects and a privacy adversary model in cryptocurrencies called the *counterparty adversary model*. On these theoretical bases, we defined an abstract model of blockchain-based cryptocurrency transaction protocols called the *coin transfer system*, and unlinkability over it called *coin transfer unlinkability (CT-unlinkability)*. Furthermore, we introduced zero-knowledgeness for the coin transfer systems to propose a method to easily prove the CT-unlinkability of cryptocurrency transaction protocols.

Finally, we proved that Zerocash is CT-unlinkable by using our proving method to demonstrate its effectiveness. It is also possible to utilize our proving method to design brand-new prospective CT-unlinkable anonymous cryptocurrencies in the future.

## Acknowledgements

## Author Contributions

All authors contributed equally to this work.

## Conflict of Interest

All authors have affirmed they have no conflicts of interest as described in Ledger's Conflict of Interest Policy.

## Notes and References

[1] Bonneau, J., Miller, A., Clark, J., Narayanan, A., Kroll, J. A., Felten, E. W. "SoK: Research Perspectives and Challenges for Bitcoin and Cryptocurrencies." In *36th IEEE Symposium on Security and Privacy (S&P)*. **1** 104–121 (2015) `http://doi.org/10.1109/SP.2015.14`.

[2] Amarasinghe, N., Boyen, X., Mckague, M. "A Survey of Anonymity of Cryptocurrencies." In *Australasian Computer Science Week Multiconference (ACSW)*. **1** 1–10 (2019) `https://doi.org/10.1145/3290688.3290693`.

[3] Silveira, A., Betarte, G., Cristia, M., Luna, C. "A Formal Analysis of the Mimblewimble Cryptocurrency Protocol." *Sensors* **21.17** 5951 (2021) `https://doi.org/10.3390/s21175951`.

[4] Jedusor, T. E. "Mimblewimble." (2016) (accessed 19 August 2022) `https://docs.beam.mw/Mimblewimble.pdf`.

[5] Poelstra, A. "Mimblewimble." (2016) (accessed 22 July 2022) `https://scalingbitcoin.org/papers/mimblewimble.pdf`.

[6] National Institute of Standards and Technology (NIST). "Guide to Protecting the Confidentiality of Personally Identifiable Information (PII)." (2010) (accessed 22 July 2022) `https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-122.pdf`.

[7] Beam community. "Beam: The Scalable Confidential Cryptocurrency." (2020) (accessed 22 July 2022) `https://docs.beam.mw/BEAM_Position_Paper_0.3.pdf`.

[8] Pfitzmann, A., Hansen, M. "A Terminology for Talking About pPrivacy by Data Minimization: Anonymity, Unlinkability, Undetectability, Unobservability, Pseudonymity, and Identity Management." (2010) (accessed 22 July 2022) `http://www.maroki.de/pub/dphistory/2010_Anon_Terminology_v0.34.pdf`.

[9] Backes, M., Kate, A., Manoharan, P., Meiser, S., Mohammadi, E. "AnoA: A Framework for Analyzing Anonymous Communication Protocols." In *26th IEEE Computer Security Foundations Symposium (CSF)*. **1** 163–178 (2013) `https://doi.org/10.1109/CSF.2013.18`.

[10] Androulaki, E., Karame, G. O., Roeschlin, M., Scherer, T., Capkun, S. "Evaluating User Privacy in Bitcoin." In *17th International Conference, Financial Cryptography and Data Security (FC)*. **LNCS 7859** 34–51 (2013) `https://doi.org/10.1007/978-3-642-39884-1_4`.

[11] Nakamoto, S. "Bitcoin: A Peer-to-Peer Electronic Cash System." (2008) (accessed 22 July 2022) `https://bitcoin.org/bitcoin.pdf`.

[12] Ben-sasson, E., *et al.* "Zerocash: Decentralized Anonymous Payments from Bitcoin." In *35th IEEE Symposium on Security and Privacy (S&P)*. **1** 459–474 (2014) `https://doi.org/10.1109/SP.2014.36`.

[13] Maxwell, G. "Confidential Transactions." (2015) (accessed 22 July 2022) `https://www.weusecoins.com/confidential-transactions/`.

[14] Parno, B., Howell, J., Gentry, C., Raykova, M. "Pinocchio: Nearly Practical Verifiable Computation." In *34th IEEE Symposium on Security and Privacy (S&P)*. **1** 238–252 (2013) `https://doi.org/10.1109/SP.2013.47`.

[15] Ben-Sasson, E., Chiesa, A., Tromer, E., Virza, M. "Succinct Non-Interactive Zero Knowledge for a von Neumann Architecture." In *23rd USENIX Security Symposium*. **1** 781–796 (2014) `https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/ben-sasson`.

[16] Groth, J. "On the Size of Pairing-Based Non-interactive Arguments." In *35th Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*. **LNCS 9666** 305–326 (2016) `https://doi.org/10.1007/978-3-662-49896-5_11`.

[17] Goldreich, O. *Foundations of Cryptography: Volume 1, Basic Tools*. Cambridge: Cambridge University Press (2001).

[18] Grin community. "Introduction to Mimblewimble and Grin." (2017) (accessed 22 July 2022) `https://github.com/mimblewimble/grin/blob/master/doc/intro.md`.

[19] Ruffing, T., Moreno-Sanchez, P. "ValueShuffle: Mixing Confidential Transactions for Comprehensive Transaction Privacy in Bitcoin." 133–154 (2017) `https://doi.org/10.1007/978-3-319-70278-0_8`.

[20] Zcash community "Zcash github." (2019) (accessed 22 July 2022) `https://github.com/zcash/zcash`.