

# Private Proof-of-Stake Blockchains using Differentially-Private Stake Distortion

Chenghong Wang, David Pujol, Kartik Nayak, and Ashwin Machanavajjhala, *Duke University* 

https://www.usenix.org/conference/usenixsecurity23/presentation/wang-chenghong

# This paper is included in the Proceedings of the 32nd USENIX Security Symposium.

August 9-11, 2023 • Anaheim, CA, USA

978-1-939133-37-3

Open access to the Proceedings of the 32nd USENIX Security Symposium is sponsored by USENIX.

# Private Proof-of-Stake Blockchains using Differentially-Private Stake **Distortion**

David Pujol Ashwin Machanavajjhala Chenghong Wang Kartik Nayak Duke University Duke University Duke University Duke University

## Abstract

Safety, liveness, and privacy are three critical properties for any private proof-of-stake (PoS) blockchain. However, prior work (SP'21) has shown that to obtain safety and liveness, a PoS blockchain must in theory forgo privacy. Specifically, to ensure safety and liveness, PoS blockchains elect parties based on stake proportion, potentially exposing a party's stake even with private transaction processing. In this work, we make two key contributions. First, we present the first stake inference attack applicable to both deterministic and randomized PoS with exponentially less running time in comparison with SOTA designs. Second, we use differentially private stake distortion to achieve privacy in PoS blockchains, and design two stake distortion mechanisms that any PoS protocol can use. We further evaluate our proposed methods using Ethereum 2.0, a widely-recognized PoS blockchain in operation. Results demonstrate effective stake inference risk mitigation, reasonable privacy, and preservation of essential safety and liveness properties.

# Introduction

Nakamoto introduced Bitcoin's protocol [47], promoting open, permissionless participation in distributed ledger maintenance. However, achieving consensus in the permissionless setting is hard as the traditional honest majority (or two-thirds majority) assumption no longer hold as the adversary can create any number of Sybil parties [23] Permissionless protocols address safety against Sybils by leveraging constraint resources, assuming honest parties hold the majority of such resources. They are classified into Proof-of-Work (PoW) protocols [47, 51, 54], reliant on computing power, and Proof-of-Stake (PoS) protocols [17, 58, 27, 42, 9, 8, 10], using monetary stake. While both protocols mark a major breakthrough in distributed consensus in a permissionless context, the reliance on expensive computing power leads to significant energy wastage in PoW [18], prompting a trend of shifting from PoW to PoS blockchains [51, 27].

The open setup of the blockchain makes transaction privacy an important issue. Several PoW designs [51], [49], [13] have taken the first step toward addressing this by concealing transaction details using cryptographic primitives while maintaining blockchain functionality. However, in PoS blockchains, merely hiding transactions is insufficient, as the stake is also sensitive [43]. Although some efforts [39, 29, 6] address PoS privacy by hiding both transaction and stake, a recent study [43] theoretically shows that SOTA designs remain susceptible to side-channel leakages, i.e., the frequency a party adding a new block, which correlates with the party's stake. Additionally, they introduce reverse tagging attacks (RTA), enabling attackers to determine if a party's stake exceeds a specified value. The basic attack primitive has been deemed feasible in real-world blockchains like Zcash [51]. Kohlweiss et al. [43] also propose a proof-of-concept stake inference attack (SIA) that estimates a party's unknown stake by repeatedly using RTA with varying comparison values.

This paper tackles the following open questions. The stake inference attack in [43] is impractical due to simplifying assumptions like deterministic protocols and strict liveness guarantees. These assumptions do not hold in real-world executions as PoS blockchains typically operate under a probabilistic manner [17, 27]. Hence, the first open question is whether there are stake inference attacks that can be launched in a real-world context. Moreover, we are not aware of any PoS design that can ensure privacy against attacks like RTA, or the derived SIA. Thus, the second open question is if there exist private PoS designs that address SIA risks. Motivated by the open questions, we elaborate on our contributions as follows:

#### 1.1 Contributions

# C-1. Stake inference attacks (SIA) for PoS protocols.

• A general stake inference paradigm. We propose a general stake inference paradigm that provides a template

for innovating efficient and accurate SIA that can be run in practice. Our paradigm captures randomized protocols by formulating liveness in a probabilistic manner and considers robust inference strategies that tolerate unreliable feedback, i.e. incorrect comparisons between a candidate value and the target stake.

• RdBin, a concrete stake inference attack. Following our inference paradigm, we provide RdBin, the first practical SIA against randomized PoS protocols. RdBin combines random walk and binary search to adaptively approximate the target stake with exponentially more efficient running time than the SOTA one [43], i.e. sublinear complexity related to the total stake S in contrast to the linear time ensured by [43].

## C-2. Private PoS using DP Stake Distortion.

- Formulating private PoS security model. In the presence of SIAs, existing PoS protocols under the standard UC security model [39, 43] are clearly inadequate to achieve privacy. To address this, we define two desiderata that ensures a practical PoS blockchain has rigorous and well-defined privacy. Based on this, we take the first attempt to integrate differential privacy (DP) [26] with previous UC security models of private ledgers and derive a formal definition for the private PoS ledger. The improved model inherits key components of UC private ledger definitions but also formulates provable stake and transaction privacy.
- · Two stake distortion mechanisms that extend standard PoS to private PoS while balancing liveliness, safety and privacy. Though naïve methods exist, i.e. partition (disconnect) each party or use equal chance for leader election, that address the potential privacy violations (SIA risks). The former completely abandons liveness while the latter provides no safety. In this work, we propose two DP stake distortion mechanisms, namely Timer and Binary mechanism, which can be mounted to any PoS protocol and extend the ordinary PoS into a private PoS under our security model. Moreover, the resulting protocols ensure provable stake and transaction privacy while preserving the original safety and liveness guarantees at the same time.
- Publicly verifiable protocols that securely realize DP distortion mechanisms. In addition to the idealized functionalities, we construct concrete protocols that securely realize the them. To the best of our knowledge, the protocols are the first attempts to provide publicly verifiable DP mechanisms over continual observations.
- Case study with real-world blockchains We conduct case studies on Ethereum 2.0 [27] to show that the stake distortion imbues real-world blockchains with reasonable stake privacy while retaining the original safety and liveness. We also compare the methods and examine how they differ in ensuring these properties.

• Prototype protocol and efficiency evaluation. We implement prototype stake distortion protocols and evaluate their efficiencies. The result indicates evaluating stake distortions adds small overheads to blockchain protocols (i.e. less than 13% of Zcash [51]'s protocol). Moreover, the distortion occurs over certain intervals, which causes the amortized cost to be even smaller.

The paper is structured as follows: Section 2 covers background, Section 3 addresses stake inference attacks, Sections 4 and 5 present private PoS desiderata and our privacy model, while Sections 6.7, and 8 detail our private PoS design, case studies, and evaluations.

# **Background**

**Private ledger.** We now examine the private ledger functionality ([2.1]),  $\mathcal{G}_{PL}$ , considered by this work, which is adapted from earlier UC formulations [39, 43] of private ledgers. Given its complexity, we outline the functionality concisely, suitable for technical dialogues. For the extended discussions, please refer to our full version [57].

## Functionality 2.1: $G_{PL}$

The functionality  $G_{PL}$  manages a general ledger state, state; a local state state; (a prefix of state) for every party  $P_i$ ; a buffer of unconfirmed transactions, buffer; and a sequence of honest inputs  $I_H$ .

Upon receiving any input I from party  $P_i$ , record  $I_H \leftarrow I_H || (I, P_i, t)$ , if I is not SUBMIT command, then evaluate the following:

- (a) (Add Transaction). If I is (SUBMIT, sid, tx), and tx is valid. Add tx to buffer, and update  $I_H =$  $I_H||(I, \mathsf{blind}(\mathsf{tx}), t)|$  with the blinded transaction, i.e. it hides the sender, recipient, and amount of tx
- (b) (Read State). Return the blinded version of  $P_i$ 's local state, blind(state<sub>i</sub>), with all non- $P_i$  transactions replaced by blind ones. If  $P_i$  is corrupted, send (blind(state<sub>i</sub>),  $I_H$ , Lkg) to the adversary.
- (c) (Extend chain). If I is (MAINTAIN\_LEDGER, sid), then perform the *ledger maintenance* [40, 19, 3] to add valid transactions in buffer to state.

In line with prior formulations [39, 43], we let the adversary know when corrupted party transactions are confirmed. Thus, in blind(state<sub>i</sub>), only non- $P_i$  transactions are replaced with blind ones. In addition,  $\mathcal{G}_{PL}$  is also parameterized with a leakage profile Lkg, revealing block proposer (the party who maintains the ledger) identities per round. It is noteworthy that Kohlweiss et al. [43] theoretically demonstrated the nonexistence of private ledgers without leaking block proposer identities (even in presence of anonymous broadcast channels). We stress that Lkg is indeed protocol-specific which especially depends on how the protocol elects the block proposer. However, in PoS blockchains, block proposers are typically chosen through a "private lottery" [29], where the odds of a party winning ledger maintenance eligibility is proportional to their stake. Hence, one may nevertheless formulate Lkg as a function over each party's stake.

**Definition 2.1.** Protocol  $\prod$  is said to be a private ledger with leakage Lkg if it UC-realizes [15]  $G_{PL}$  in the presence of a p.p.t. adversary.

Private ledger properties. Following the ledger properties discussed in previous works [40, 19, 39, 43], we formulate important properties considered by this work.

**Definition 2.2** (Private ledger properties). Let  $\prod$  to be a private ledger protocol executed by a set of parties. We consider the following properties hold when f fraction of the parties is malicious.

- Safety: No two honest users disagree on confirmed transactions (i.e., transactions appended to state).
- (z,t)-Liveness: Any valid transaction that is input to at least an honest  $z \le (1-f)$  fraction of the total parties, will be appended to State within t time steps.

Consistent with [17] [14] [45] [59], we consider at least  $\frac{2}{3}$  fraction of the total stake is owned by the honest parties, and correspondingly  $f < \frac{1}{3}$  is interpreted as malicious corrupted stake fraction. Moreover, for various reasons, the properties in Definition [2.2] may be ensured in a probabilistic manner. For example, Algorand [17] and Ouroboros [40] use randomness for efficiency thus ensures safety or liveness only with high probability.

Privacy tradeoffs and attacks against PoS. Kohlweiss et al. [43] demonstrated that it is theoretically impossible to have a private ledger ensuring liveness without revealing the block proposer's identity. Moreover, recognizing block proposers over time reveals the frequency at which parties take on the block proposer role. In PoS blockchains, such frequency correlates with a party's stake, potentially leading to stake privacy loss and allowing attackers to deduce an individual's stake. For example, Kohlweiss et al. [43] examined the Reverse Tagging Attack (RTA), enabling attackers to assess if a party's stake surpasses a specified threshold. RTA's idea involves creating input disparities between a target and other nodes, such as by delaying specific transactions. For example, the attacker delays a self-created transaction tx for a target  $P_v$  with stake  $f_v$ , and a set of corrupted nodes with stake  $f_{cmp}$ . If tx isn't confirmed after t steps, by (z,t)-liveness, the attacker learns  $f_v \ge 1-z-f_{\text{cmp}}$ .

The delay operation was implemented using the *In-vblock* technique, first proposed in [46] and was later used to infer Bitcoin topology [21]. Kohlweiss et al. [43]

further assessed the validity of *Invblock* against private ledgers, i.e. Zcash [51]. Note that, the cost of launching *Invblock* is low, as it necessitates the attacker becoming one of the victim's peers in the P2P network, rather than compromising all target node peers.

Additionally, [43] also outlined a stake inference attack (SIA) design by iteratively applying RTA with adaptive  $f_{cmp}$ . The attacker starts with  $f_{cmp} = 0$  and repeats RTA, increasing  $f_{cmp}$  by  $\frac{1}{S}$  each time until the first  $f_{cmp}$  satisfies  $f_{v} \ge 1-z-f_{cmp}$ , where  $\hat{f}_{v} = 1-z-f_{cmp}$  approximates  $f_{v}$ .

**Differential privacy** [26]. guarantees that altering one input data element in an algorithm or mechanism imposes minimal impact on the output. More specifically, let D and D' be two databases that differ by only a single tuple, namely neighboring databases, then:

**Definition 2.3** ( $\epsilon$ -differential privacy). Given  $\epsilon > 0$ , a randomized mechanism  $\mathcal{M}$  is  $\epsilon$ -DP if for all pair of neighboring databases D,D', and any possible output  $o \subset Range(\mathcal{M})$ , the following holds:

$$\Pr[\mathcal{M}(D) \in o] \le e^{\epsilon} \Pr[\mathcal{M}(D') \in o]$$

An alternative definition, known as *approximate DP* or  $(\epsilon, \delta)$ -DP, allows a small failure probability  $\delta$  in addition to the constraints in Definition 2.3 In this paper, the focus is on  $\epsilon$ -DP, while an extended discussion that considers the relaxed  $(\epsilon, \delta)$ -DP can be found in our full version 57.

## 3 Stake Inference Attack

In this section, we discuss practical stake inference attacks against PoS protocols. In general, given a target party  $P_{\nu}$  with stake  $f_{\nu}$ , we consider attack strategies to adaptively choose  $\hat{f}_{\nu}$  that minimizes the error  $|\hat{f}_{\nu} - f_{\nu}|$ .

# 3.1 Issues with the SOTA Approach [43]

- **I-1. Restricted to deterministic protocols.** Kohlweiss et al. [43] consider rather idealized deterministic protocols, assuming stringent liveness conditions. However, this assumption may not align with real-world scenarios, as most blockchains function probabilistically, with liveness typically guaranteed at a high probability [17] [27].
- **I-2. RTA: One-sided comparison only.** Based on (z,t)-liveness, tx can be confirmed even if less than z fraction of honest parties receive it. Therefore, the  $s_v + f_{cmp} < (1-z)$  conclusion may not hold when tx is confirmed within t steps, potentially causing  $f_v$  to be larger than  $\hat{f}_v$ . This one-sided comparison limitation makes RTA-based SIA less reliable for accurate stake inferences.
- **I-3. Inefficient search strategy.** The linear scan strategy considered by [43] requires a large number of RTA executions, which is inherently inefficient.

The aforementioned limitations reduce the practicality of RTA-based SIA. In what follows, we focus on developing a practical SIA that is efficient, accurate, and suitable for randomized PoS protocols.

# The General Stake Inference Paradigm

In this section, we show that the problem of finding accurate stake inference against probabilistic PoS blockchains can be reduced to a variant of *Noisy Search* problem [37].

**Definition 3.1** (Noisy search problem). Given n coins sorted by head probabilities ( $p_i$  for the  $i^{th}$  coin) and a target coin with  $p_1 \le p^* \le p_n$ , an algorithm must find two coins i, i+1 satisfying a given  $\tau > 0$  such that  $[p_i, p_{i+1}]$  intersects  $[p^*-\tau, p^*+\tau]$ , without knowing exact head probabilities but being allowed to flip coins

The general reduction idea is to map stake to simulated biased coins so that inferring stake values reduces to locating search coins that is close to a target coin with heads probability tied to  $f_{\nu}$ . Details are as follows.

Formulating liveness. To simulate stake-based biased coins, we rely on the liveness ensured by the blockchain. To apply to randomized protocols, we adapt the (z,t)-Liveness definition accordingly.

**Definition 3.2** (Probabilistic (z,t)-liveness). For any transaction tx that is input to z fraction of the honest parties, the probability tx is NOT confirmed after t time steps equals to q(1-z,t), where  $q(\cdot)$  is a monotonically increasing function related to 1-z.

In general,  $q(\cdot)$  is highly related to the protocol specification. For example, assuming the probability of a party proposing the next block is determined by a slot function  $\phi(\cdot)$  [29, 39, 43] over its stake. One may obtain  $q(1-z,t) = (1-\phi(z))^t$  for any z,t > 0. Although, this probability may differ from the theoretical one in real executions, i.e., due to some other configurations, such as timeouts, transaction pool size, etc. One may assume that as long as the protocol specification does not change, the value of q(z,t) remains stable for any z and t. Moreover, the actual value of  $q(\cdot)$  in real executions may be not directly computable, but one may evaluate it empirically.

Simulating "stake-based" biased coins. Given a node (or a group of nodes) with relative stake  $f_x$ , we provide an interface for simulating a biased coin with heads probability equal to  $q(f_x,t)$  for any t>0. Note there are two types of biased coins, the search coin and the target coin, for which the simulations are different. In simulating the search coin, the attacker knows the pre-image  $f_x$ , while it is unknown when simulating the target coin.

To simulate a search coin, the attacker first samples a set of corrupted nodes with total stake  $f_x$ , then broadcasts a self-created transaction, tx, to everyone, but removes tx from the sampled nodes' transaction pool. The attacker waits for t time steps to observe if tx is confirmed. By Definition 3.2 the probability of tx not being confirmed after t steps is  $q(f_x,t)$ . For the target coin, we utilize the delay operation considered by previous works [43, 46]. Specifically, the attacker broadcasts tx to everyone but delays it to the victim, such that tx does not reach  $P_v$  for at least t time steps. Similarly, the attacker then waits to observe if tx is confirmed after t steps, where the "not confirmed" probability equals  $q(f_v,t)$ .

For demonstration purposes, we consider t = 1 to be the default setting and therefore omit it. In what follows, we use  $p_x$  to directly denote the heads probability of a stake-based biased coin with pre-image  $f_x$ .

**Reduction.** We provide an abstract algorithm for SIA by assuming the existence of a black-box solver, ns\_solver, for the noisy search problem. Given  $P_{\nu}$  with stake  $f_{\nu}$ , and n sorted stake values, the algorithm simulates corresponding search coins as well as the target coin then inputs them to ns\_solver along with a specified  $\tau$ . Upon receiving two coins i, i+1 output from ns\_solver, the algorithm computes  $q^{-1}(p_i)$  and,  $q^{-1}(p_{i+1})$  to obtain the pre-images  $f_a$ , and  $f_b$ , then returns  $\frac{f_a+f_b}{2}$  as  $\hat{f}$ . Note that, in simulating the search coins, the algorithm knows the pre-image of each stake-based biased coin, and thus  $\frac{f_a+f_b}{2}$  can be obtained within constant time. Therefore, if  $ns\_solver$  requires T(n) times coins tosses before termination, then the running time of the SIA is dominated by T(n) many simulated coin flips. Figure 1 shows a diagram of the aforementioned reduction.

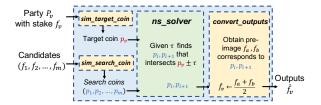


Figure 1: Reduction diagram.

#### 3.3 **Practical Stake Inference Attack**

Following our reduction, we provide RdBin, a practical SIA that combines random walk and binary search. Similar to [43], we examine a loosely dynamic scenario where honest parties' stakes remain stable over short periods but change over longer durations. RdBin is parametered by  $\theta > 0$ , and  $\tau \in (0,1)$ , and is given a balanced (stake) binary tree ST with each node labeled with a stake segment  $[f_a, f_b] \subseteq [0, f]$ . Specifically, the root is labeled with [0, f], and for every internal node v with label  $[f_a, f_b]$ , its left and right child is labeled with  $[f_a, m]$  and  $[m, f_b]$ , respectively, where  $m = \lfloor \frac{f_a + f_b}{2} \rfloor$ . Moreover, all leaf nodes satisfy  $f_b - f_a = \frac{\theta}{S}$ . The algorithm starts with a random walk in ST from the root, and for each round t (assuming at node v(t)), it evaluates a probabilistic stake comparator, pSC, that compares  $f_v$  with the pivot of v(t), say  $\frac{f_a+f_b}{2}$ . If  $f_v < \frac{f_a+f_b}{2}$  is claimed by pSC, it moves to the left child, otherwise to the right child. RdBin halts and outputs  $\hat{f} \leftarrow \frac{f_a+f_b}{2}$ , if it reaches a leaf node or pSC asserts min\_diff. The details are shown in Algorithm [1].

# **Algorithm 1** pSC $(f_v, f_{cmp}, \tau, \delta)$

```
1: Sample a set of corrupted users P_c with stake f_{cmp}

2: for i=1,2,3... do p_0=p_1=0

3: \tau_i=e^{-i/2}, \, \delta_i=\delta/e^i, \, n_i=\log(1/\delta_i)/\tau_i^2

4: if \tau_i<\tau then assert min_diff

5: for j=1,...,n_i do (b\in\{0,1\})

6: Broadcast conflicting transactions ^{\dagger} tx<sub>0</sub>, tx<sub>1</sub>

7: Delay tx<sub>0</sub> to P_v and remove tx<sub>1</sub> from P_c.

8: p_b=p_b+1/n_i, if tx<sub>b</sub> is not confirmed

9: if |p_0-p_1| \geq 2\tau_i then

10: assert f_v>f_{cmp} if p_0>p_1 else f_v\leq f_{cmp}
```

<sup>†</sup>As long as one transaction is confirmed, the other one is invalid, this allow us to flip two biased coins at the same time.

**Theorem 1.** Given  $\delta$ ,  $\tau \in (0,1)$ ,  $\tau_c = |q(f_v) - q(f_{cmp})|$ , and  $\tau_m = \max(\tau, \tau_c)$ , with probability at least  $1 - \delta$ , pSC outputs correctly after expected  $O\left(\tau_m^{-2}\ln(1/\delta\tau_m)\right)$  many simulated coin tosses.

For each round i, the failure probability of testing coins is at most  $\delta_i$ , thus by union bound [32], the failure probability for the entire testing is bounded by  $\sum_i \delta_i < \delta$ . By Hoeffding's inequality, the probability that pSC to continue running for  $\tau_i < \frac{1}{2}\tau$  decreases exponentially in  $n_i$ . Thus the expected complexity is dominated by the complexity of round  $i = \log(\frac{1}{\max(\tau,\tau_c)})$ , which is consistent with Theorem [1] We extend this to full proofs in Appendix [B.1]. Overall, pSC provides accurate two-sided comparisons with a bounded failure probability. Moreover, by setting  $\delta = (\log \frac{S}{\theta})^{-1}$ , the random walk ensures movement towards the correct node with probability at least  $1-(\log \frac{S}{\theta})^{-1}$ , and by union bound, the failure probability of RdBin is bounded by a constant factor.

**Theorem 2.** Let  $n = \frac{S}{\theta}$ ,  $\delta = O(1/\log n)$ ,  $p_v = q(f_v)$ , and  $\eta = \max(|q^{-1}(p_v \pm \tau) - f_v|)$ , then RdBin satisfies:

- 1. The running time is  $O(\tau^{-2}\log(\log n/\tau)\log n)$ .
- 2. With constant failure probability, RdBin outputs an inference  $\hat{f}_v$  with error  $|\hat{f}_v f_v| \le \max(\frac{\theta}{S}, \eta)$ .

Note that RdBin represents only one construction under our general inference paradigm, and there exist many other designs in this space. Due to the simplicity of RdBin and the fact that its implementation does not require extra memory, we present it as the default attack. One may also derive certain extensions on RdBin. First, it is possible to flip even more simulated coins at once. For instance, one can create a set of (conflicting) transactions,

 $\mathsf{tx}_1, \mathsf{tx}_2, ..., \mathsf{tx}_n$ , broadcasting them all at once, but delaying the arrival of each  $\mathsf{tx}_j$  to a specific target  $P_j$ . This further enables the attacker to construct a parallelized RdBin against multiple victims. Second, as the stake is the additive outcome of a series of transactions. It's obvious that one can utilize RdBin to extrapolate transactions outcome for any victim  $P_v$  between any two time steps.

#### 4 Private PoS Desiderata

We discuss requirements for achieving a private PoS by defining two desiderata that are designed to ensure a practical protocol has rigorous and well-defined privacy.

**D-1:** Well-defined privacy on leakage. Any private PoS should provide a rigorous and provable bound on the leakage of each party's stake and transactions. To facilitate this we say that in addition to the standard UC security model (Definition 2.1), a private PoS should also meet additional privacy definitions, i.e., Definition 5.4 which constrains the knowledge that an adversary can obtain regarding each party's stake and transactions by observing the associated leakages.

**D-2. Balance privacy, safety and liveness.** While safety and liveness are traditionally ensured through a PoS protocol, introducing privacy may create tension between these properties. For instance, one could ensure privacy by using a consensus protocol that is completely independent of the parties' stake. But such a system would not be safe, since the fundamental purpose of introducing stake is to address safety against Sybils. Thus, given all of safety, liveness, and privacy may not be achievable simultaneously, we require a practical private PoS ledger should balance between these three properties.

# 5 Privacy Model

## 5.1 Privacy Requirements

**PR-1.** Limited stake inference. In general, we require that any attacker should not infer the current stake owned by a certain party (by observing the protocol leakage) within an additive bound  $\alpha$ . We assume that the adversary knows all parties involved in the protocol, as their participation can be easily identified, for example, when block proposers broadcast the updated ledger states, it inevitably indicates their involvement in the ledger maintenance. Hence, we do not require stronger privacy in the form of "the existence of a party cannot be inferred".

**PR-2. Privacy with expiration.** Considering that the parties' stake changes dynamically, stake privacy should be ensured under continual observation [25]. Ideally, we desire to obtain the same degree of privacy for each time step based on all information disclosed until that time step. However, providing such strong privacy guarantees

inherently leads to errors accumulating over time [25, 24]. In the context of PoS blockchains, this can result in a large distorted adversarial stake potentially causing a significant decline in safety or liveness guarantees. Therefore, to balance between these properties, in this work, we adopt a relaxed "privacy with expiration" model [12, 38]. In this model, we aim to achieve strong privacy for parties' stake within a recent time frame while allowing the privacy of parties' past stake to gradually decay over time.

PR-3. Transaction privacy. We also consider transaction privacy, that is, we bound the attacker's ability to infer the outcome of any party's transaction. Typically, as the stake value directly reflects the additive outcome of transactions, achieving PR-1 and PR-2 also implies transaction privacy if the protocol leakage is subject to parties' stake, i.e., no information relevant to the transaction outcome is disclosed to the adversary, other than what can be obtained from the parties' stake at each time.

Putting it all together. Combining PR-1, 2, 3, we formulate our privacy requirements using Pufferfish [41] framework, a well-known tool that translates common privacy requirements into formal definitions [35].

**Definition 5.1** (Stake privacy requirement). Let P = $\{P_1, P_2, ..., P_n\}$  denote the set of participants with  $S^i =$  $\{s_1^i, s_2^i, ...\}$  as the stake history for party  $P_i$ , where  $s_i^i$  is  $P_i$ 's stake at time j. We denote  $\theta_i \in \Theta$  as the attacker's belief about  $P_i$ 's stake, where  $\Theta$  is the set of all possible adversarial beliefs. Given a private ledger with leakage Lkg, and  $\epsilon, \alpha > 0$ , we require that at any time t the following holds for all  $\theta_i \in \Theta$ , and  $o \in Range(Lkg)$ .

$$\frac{\Pr_{\theta_i} \left[ s_j^i = x \mid \mathsf{Lkg} = o \right]}{\Pr_{\theta_i} \left[ s_j^i = y \mid \mathsf{Lkg} = o \right]} / \frac{\Pr_{\theta_i} \left[ s_j^i = x \right]}{\Pr_{\theta_i} \left[ s_j^i = y \right]} \le e^{\omega(t-j)\epsilon} \tag{1}$$

where  $j \in [1,t]$  and  $x,y \in \mathbb{R}$  such that  $\Pr[s_i^i = x] \ge 0$ ,  $\Pr[s_i^i = y] \ge 0$ , and  $x \le y \le x + \alpha$ .  $\omega$  is the privacy decay multiplier, which is a monotonically increasing function related to (t - j) with  $\omega(0) = 1$ .

As in Definition 5.1, the privacy requirement limits the maximum Bayes factor that an attacker can learn, after observing the leakage Lkg, regarding the stake owned by any party at any point of time. Such a factor with respect to the most recent stake  $s_t$  is bounded by  $e^{\epsilon}$ . Furthermore, for any stake  $s_i$  that is (t - j) time steps away from current time, such a factor is bounded by  $e^{\omega(t-j)\times\epsilon}$ , which captures the notion of privacy decay (PR-2). Note that, Definition 5.1 also addresses transaction privacy (PR-3). More specifically, let  $tx_{j:j+1}$  as  $P_i$ 's additive transactions outcome between time j and j + 1. Since  $tx_{i:i+1} = s_{i+1} - s_i$  and by Equation 1, the attacker's posterior odds (after observing the leakage Lkg) about  $\mathsf{tx}_{j:j+1} = x$  rather than  $\mathsf{tx}_{j:j+1} = y$   $(x \le y \le x + \alpha)$  is at most  $e^{\omega(t-j)\times\epsilon}$  times the attacker's prior odds.

#### 5.2 **Formal Privacy Definition**

In this section, we discuss the formal privacy model considered by this work. First, we formulate the leakage privacy following the notion of differential privacy [26].

**Definition 5.2** ( $(\alpha, u)$ -neighbors). Let  $S_t = \{s_1, ..., s_t\}$ and  $S'_t = \{s'_1, ..., s'_t\}$  be the any two stake histories up to time t.  $S_t$  and  $S'_t$  are  $(\alpha, u)$ -neighbors if the following holds: (i)  $\alpha > 0$ , and  $u \in (0,t]$ ; (ii)  $s_j = s_j'$  for all  $j \le u$ ; and (iii)  $s_j \le s'_i \le s_j + \alpha$  for all j > u.

**Definition 5.3** ( $(\alpha, \epsilon)$ -private leakage). Given Lkg that depends on stake, Lkg is said to be  $(\alpha, \epsilon)$ -private, if for any two  $(\alpha, u)$ -neighbor stake histories  $S_t$ , and  $S'_t$ , any output  $o \subset Range(Lkg)$ , and any  $P_i$ , the following holds

$$\Pr\left[\mathsf{Lkg}^{\langle P_i, S_t \rangle} = o\right] \le e^{\omega(t-u)\epsilon} \cdot \Pr\left[\mathsf{Lkg}^{\langle P_i, S_t' \rangle} = o\right] \tag{2}$$

where  $\mathsf{Lkg}^{\langle P_i, S_t \rangle}$  (resp.  $\mathsf{Lkg}^{\langle P_i, S_t' \rangle}$ ) denotes the leakage when  $P_i$ 's stake history is activated by  $S_t$  (resp.  $S'_t$ ), and  $\omega$  is the decay function related to t - u with  $\omega(0) = 1$ .

**Theorem 3.** If Lkg is  $(\alpha, \epsilon)$  private (Definition 5.3), then Lkg satisfies all privacy requirements (Definition 5.1).

We stress that DP is a special case of Pufferfish privacy [41], and applying DP on party's stake history ensures all privacy requirements listed in Section 5.1. For complete proofs, interested readers can refer to Appendix B.3. Next, we formally define the private PoS ledger under the UC framework [15].

**Definition 5.4**  $((\alpha, \epsilon)$ -private PoS ledger). A PoS ledger protocol is said to be  $(\alpha, \epsilon)$ -private (resp.  $(\alpha, \epsilon, \delta)$ private), if: (i) it UC realizes [5]  $G_{PL}$  with leakage Lkg in the presence of a p.p.t. adversary A and (ii) the leakage Lkg is  $(\alpha, \epsilon)$ -private (resp.  $(\alpha, \epsilon, \delta)$ -private).

We stress that private ledger protocols are allowed to interact with auxiliary functionalities (hybrids), which captured the resources that are available to parties. The composability property [5] of UC states that if  $\pi$  interacts with  $\mathcal{F}$  to UC realize functionality  $\mathcal{G}$ , and  $\pi_{\mathcal{F}}$  UC realizes  $\mathcal{F}$ , then substituting calls to  $\mathcal{F}$  in  $\pi$  with calls to  $\pi_{\mathcal{F}}$ results in a secure protocol for G in the hybrid world. This also provides the flexibility to analyze the security of complex protocols in a modular manner. In this paper, we primarily focus on the ledger maintenance protocol (i.e., the sub-protocol that UC realizes MANTAIN\_LEDGER command in  $G_{Pl}$ ), as this is the main component that leaks stake-related information [19, 29, 39].

## 6 Private PoS with Stake Distortion

To mitigate SIA risks, naïve approaches exist. For instance, one may consider each party to maintain the ledger with the same chance, but clearly, this provides no safety at all. Although this can be improved by letting one party own multiple nodes, each contributing independently to ledger maintenance, it requires a single party running the private lottery protocol multiple times, causing significant overhead for large stakeholders. Another option is to disconnect or impose long delays for each party to produce new blocks, but this severely impacts liveness. Steered by our privacy model, we propose our private PoS design, which not only ensures provable privacy but also strikes a balance between safety and liveness.

# 6.1 Design Overview

We first provide a brief overview of the general design pattern (Protocol 6.1) for a PoS-based ledger maintenance protocol, from the view of a party P. This pattern is abstracted from previous constructions [43] [29].

# **Protocol 6.1: Ledger Maintenance,** $\prod_{LM}$

- Get lottery inputs: (in, com<sub>in</sub>, r) ← F<sub>stk</sub>
  // in obtained from F<sub>sd</sub> is the true stake
  (In our design, ∏<sub>LM\*</sub>, lottery inputs are obtained via interacting with (in, com<sub>in</sub>, r) ← F<sub>sd</sub>, where in returned by F<sub>sd</sub> is the distorted stake)
- 2: Evaluate a private lottery: ret ← priv\_lottery(in)
- 3: if ret == WIN then
- 4: Generate zero-knowledge proof,  $\pi$ , proving:
  - (a) the party wins the lottery with input in
  - (b) lottery input in is consistent with com<sub>in</sub>.
- 5: Perform the ledger maintenance activities.

In general, ledger maintenance execution occurs in discrete time slots, during which slot leaders are elected through private lotteries to update ledger states. Specifically, in each round, each party P interacts with a functionality  $\mathcal{F}_{\text{stk}}$  [29] to obtain the lottery input in and partakes in a private lottery with winning odds tied to in. If P wins, then she contributes to ledger maintenance along with a zero-knowledge proof of eligibility. Typically, to implement such a proof, other parties should also access a hiding commitment  $com_{in}$  to P's lottery input from  $\mathcal{F}_{\text{stk}}$  (while in practice, such commitments can efficiently be computed from the blockchain [39]). However, the actual value of in and the opening r for  $com_{in}$  remain exclusive to P. In a PoS blockchain, the lottery input in returned by  $\mathcal{F}_{\text{stk}}$  is P's **true stake** at the current time.

However, as mentioned before, the standard PoS design unavoidably reveals the frequency of each party's lottery wins, closely linked to the party's true stake and enabling accurate stake inferences. To mitigate this, our design modifies how parties acquire lottery inputs while

imposing no changes to other components. Specifically, we consider parties interacting with a new functionality,  $\mathcal{F}_{sd}$ , which produces a **distorted stake** as the lottery input for each party at every time. This suggests that one may directly derive an actual protocol under our design from the standard one by mounting functionality  $\mathcal{F}_{sd}$ , i.e., replacing how the parties obtain private lottery inputs with a subroutine call to  $\mathcal{F}_{sd}$ . As the input retrieval is completely independent of the core protocol, the derived protocol thus retains the security guarantees of the original protocol [29]. Moreover, the leakage for the derived protocol is no longer related to users' true stake, but rather tied to their distorted stake.

**Theorem 4.** Given protocol  $\prod_{LM}$  that UC-emulates MAINTAIN\_LEDGER in  $\mathcal{G}_{PL}$  with  $\mathcal{F}_{stk}$  and  $\mathsf{Lkg} = \phi\left(\{S^i\}_{1 \leq i \leq n}\right)$ , replacing  $\mathcal{F}_{stk}$  in  $\prod_{LM}$  with  $\mathcal{F}_{sd}$  yields a protocol  $\prod_{LM^*}$  that UC-emulates MAINTAIN\_LEDGER with  $\mathsf{Lkg} = \phi\left(\{\tilde{S}^i\}_{1 \leq i \leq n}\right)$ , where  $\phi$  is a slot leader function mapping parties' stake distribution over time to a sequence of slot leaders, and  $S^i$  and  $\tilde{S}^i$  denote  $P_i$ 's true and distorted (by  $\mathcal{F}_{sd}$ ) stake history, respectively.

We provide complete proof of Theorem [4] in Appendix [B.4]. As stake-related leakage occurs solely in leader election [40], [39], [19], we consider the distorted stake to be used exclusively for the private lottery. For instance, to issue and validate payment transactions, one should keep using the true stake. This further suggests that other protocols of the standard design, i.e., the transaction submission and validation, undergo no alterations. In addition, by the composability of UC, if  $\prod_{LM^*}$  realizes  $\mathcal{G}_{PL}$  through calls to  $\mathcal{F}_{sd}$ , and  $\prod_{sd}$  realizes  $\mathcal{F}_{sd}$ , then substituting calls to  $\mathcal{F}_{sd}$  with subroutine calls to  $\prod_{sd}$  leads to a secure protocol for  $\mathcal{G}_{PL}$ . This enables us to explicitly establish our design objectives as designing practical  $\mathcal{F}_{sd}$  that produces noisy stake satisfying Definition [5.3] and devising protocols that securely realize  $\mathcal{F}_{sd}$ .

# 6.2 Differentially Private Stake Distortion

We start with two different designs of  $\mathcal{F}_{sd}$ , namely, the timer  $(\mathcal{F}_{Timer})$  and the binary mechanism  $(\mathcal{F}_{Bin})$ .

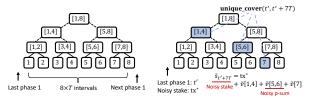
**Timer mechanism**  $\mathcal{F}_{\mathsf{Timer}}$ . In general, the mechanism periodically (every T time steps) distorts each party's stake with fresh Laplace noise drawn from  $\mathsf{Lap}(\frac{\alpha}{\epsilon})$ . For other times between two distortions, it reuses the noisy stake in the previous round. The rationale for distorting stake in a periodic manner (as opposed to distorting every time) is twofold. First, parties' stake generally exhibits stability over short intervals. Insignificant stake changes exert negligible effects on their lottery winning odds, rendering the release of a new distorted stake each time unwarranted. Second, distorting stake in a less frequent manner yields a smoother privacy decay (Section 6.3), otherwise stake privacy may drop significantly over a

short time. Note that the distorted stake can be negative due to the symmetry of the Laplace noise. We stress that a party with a negative stake is treated the same as no stake, and thus the corresponding party has no chance of winning the following private lottery.

**Theorem 5.** Protocol  $\prod_{\mathsf{LM}^*}$  mounted with  $\mathcal{F}_{\mathsf{Timer}}$  is  $(\alpha, \epsilon)$ -private with privacy decay  $\omega(t-u) = \lfloor \frac{t-u}{T} \rfloor$ .

By the theory of the Laplace mechanism [26], the privacy loss for all stake values within the latest T time steps is bounded by  $\epsilon$ . Moreover, due to multiple releases of the noisy stake at different times, the privacy loss of any historical stake follows the privacy loss under k-fold composition rules [26]. Thus, one can obtain a linear decay function such that  $\omega(t-u) = \lfloor \frac{t-u}{T} \rfloor$ . We defer the complete proof of Theorem 5 in Appendix B.5

**Binary mechanism**  $\mathcal{F}_{\mathsf{Bin}}$ . We improve upon  $\mathcal{F}_{\mathsf{Timer}}$  and provide the binary (tree) mechanism,  $\mathcal{F}_{\text{Bin}}$ , which provides a smoother (logarithmic) privacy decay over time. In general, the mechanism relies on the fact that any party's stake is the additive outcome of all her transactions. For each party  $P_i$ , the functionality internally groups all  $P_i$ 's transactions based on transaction time, wherein the additive outcome of each group represents a partial sum (p-sum) of  $P_i$ 's current stake. Subsequently,  $\mathcal{F}_{Bin}$  distorts each p-sum using Laplace noises, and for each time, the noisy stake is obtained by aggregating a set of noisy p-sums. Moreover, a single noisy p-sum may be reused to derive multiple noisy stake at different times, which further leads to less privacy loss over time.



(a) Binary interval tree

**(b)** Determine noisy p-sums

Figure 2: Example of Binary Mechanism with L = 8T

Specifically,  $\mathcal{F}_{Bin}$  features two distortion phases. In phase one,  $\mathcal{F}_{Bin}$  releases a new distorted stake every  $L \gg T$ time steps acting as a single noisy p-sum, reused for the next L steps. Phase two occurs between phase one distortions, with a new stake released every T time steps. Precisely,  $\mathcal{F}_{Bin}$  employs a binary interval tree, with leaf nodes representing time intervals of length T and internal nodes integrating their children's time intervals. Figure 2a shows an example. For every T time steps,  $\mathcal{F}_{Bin}$  identifies a set of disjoint tree nodes that *uniquely* cover the interval from the last phase one distortion to the current time. More specifically given the last phase one distortion time t' and the current time j, we use the following algorithm to cover [t, j]: (i) Set index a = t', identify the sub-interval (tree node)  $[a,b] \subseteq [a,j]$  encompassing  $T \cdot 2^k$  time steps, with  $T \cdot 2^{k+1} > |j-a|$ ; (ii) Update a = b + 1 and repeat step (i) until the entire interval is covered. For better illustration, we show an example in Figure 2b. Note that, by running the aforementioned strategy, at most  $\log_2(\frac{L}{T})$  nodes are required to uniquely any intervals. Subsequently,  $\mathcal{F}_{Bin}$  groups transactions based on selected nodes, computes noisy p-sums, and aggregates them into a new noisy stake.

For any stake  $s_u$  with  $u \le t$  (current time), the total privacy loss of  $s_u$  under  $\mathcal{F}_{Bin}$  comes from both distortion phases. By Theorem 5, the phase one privacy loss composes linearly. In phase two, the transaction outcome at time *u* impacts at most  $\log_2(\frac{\max(t-u,L)}{T})$  noisy p-sums, thus leading to a logarithmic privacy decay. Combining the two we conclude the overall privacy decay as follows.

**Theorem 6.** Protocol  $\prod_{LM^*}$  mounted with  $\mathcal{F}_{Timer}$  is  $(\alpha, \epsilon)$ -private with privacy decay

$$\omega(t-u) \le \begin{cases} \left\lfloor \frac{t-u}{L} \right\rfloor + \log_2(\frac{L}{T}), & if (t-u) > L \\ \log_2(\frac{t-u}{T}), & otherwise \end{cases}$$
(3)

#### 6.3 **Analysis of the Stake Distortion**

Resilience to SIA attacks. We first analyze how SIA errors change when stake distortion is employed. Theorem 2 implies that in standard protocols where no stake distortion is present, the SIA error can be made arbitrarily small as long as the attacker can flip the stake-based coins unlimited times over a long period of time (this allows the attackers to set  $\theta$  and  $\tau$  arbitrarily small). When switching to distorted stake, SIA only finds the interval where the distorted stake is located. At this point, the errors are based on two factors, namely, the errors inherent in the attack, and the errors due to the injected DP noises.

**Theorem 7.** Given protocol  $\prod_{LM}$  with SIA error Err. Extending it to  $\prod_{\mathsf{LM}^*}$  with  $\mathcal{F}_{\mathsf{Timer}}$  and  $\mathcal{F}_{\mathsf{Bin}}$  yields SIA errors in  $\mathsf{Err} + O_p\left(\frac{\alpha}{\epsilon}\right)^I$  and  $\mathsf{Err} + O_p\left(\frac{\alpha}{\epsilon}\sqrt{\log_2(\frac{L}{T})}\right)$ , respectively.

Theorem 7 is the direct application of Chernoff bound [36] (complete proofs in Appendix B.7) to the injected Laplace noises. Note that the SIA error bounds now rely on both the attacker-selected parameters ( $\theta$  and  $\tau$ ) and the privacy parameters ( $\alpha$  and  $\epsilon$ ), indicating that one may not achieve arbitrary precision in stake inferences with an unlimited number of simulated coin flips.

Safety analysis. It is possible that the adversary's stake after distortion, say f, is larger than the pre-distortion value f. However, as long as there is a limited number of parties (existing protocols typically bound this by imposing a minimum staking amount [17, 27], say v), the

 $<sup>{}^{1}</sup>O_{p}$  is the Big O in its probability notion [22]

bounded variance property of Laplace noise [26] enables the derivation of a slack, i.e. by Chernoff inequalities, that bounds the stake increment  $\tilde{f}-f$  with a high probability. By setting the target as  $\tilde{f}<\frac{1}{3}$ , and with the derived slack, one can obtain another threshold  $\xi$  such that as long as  $f<\xi$ , then  $\tilde{f}<\frac{1}{3}$  with high probability, i.e., the safety remains valid after stake distortion. We summarize the safety bounds in Theorem [8] and show in Section [7] that they are reasonable through a case study.

**Theorem 8.** Assuming  $0 < \beta < 1$ ,  $\Gamma = \frac{1.74\alpha}{\epsilon S} \sqrt{\frac{S}{v}} \log \frac{1}{\beta} < 1$ , and  $f = \frac{1-\gamma}{3}$  with  $\Gamma \le \gamma < 1$  (or  $\gamma \ge \sqrt{\log_2(\frac{L}{T})} \times \Gamma$ ), the adversary's stake after distortion,  $\tilde{f}$ , is less than  $\frac{1}{3}$  with a probability of at least  $1 - \beta$  under  $\mathcal{F}_{\text{Timer}}$  (or  $\mathcal{F}_{\text{Bin}}$ ).

Theorem 8 demonstrates that for a protocol,  $\prod_{LM}$ , tolerating up to  $\frac{1}{3}$  malicious stake, employing stake distortion with  $\mathcal{F}_{\mathsf{Timer}}$  and  $\mathcal{F}_{\mathsf{Bin}}$  results in a modified protocol,  $\prod_{\mathsf{LM}^*}$ , which tolerates up to  $\frac{1}{3}$ - $O_p\left(\frac{\alpha}{\epsilon S}\sqrt{\frac{S}{\nu}}\right)$ , and  $\frac{1}{3}$ - $O_p\left(\frac{\alpha}{\epsilon S}\sqrt{\frac{S}{\nu}}\log_2(\frac{L}{T})\right)$  malicious controlled stake, respectively (complete proofs in Appendix 8.8). Moreover, Theorem 8 highlights a safety-privacy trade-off, where weaker privacy guarantees (small  $\epsilon$  or large  $\alpha$ ) result in smaller  $\gamma$ , allowing for greater malicious tolerance. This offers practitioners the flexibility to balance safety and privacy, opting to exchange some safety for enhanced privacy, or vice versa.

**Liveness analysis.** Typically, the liveness guarantee of PoS ledgers is directly tied with the honest majority assumption (same as the safety guarantee) if no changes are made to the networking layer. For instance, by 19 Theorem 9 and 17 Theorem 4 the liveness is ensured for PoS protocols as long as the honest parties control more than a certain fraction of the total stake (i.e.,  $> \frac{2}{3}$ ). Note that our stake distortion does not alter the network assumptions (i.e., a network of authenticated multicast channels with bounded delayed 29, and by Theorem 1, the honest majority assumption still holds after stake distortions. Therefore, without loss of generality, in what follows we focus on the safety analysis and assume that once safety is preserved, liveness is also retained.

Impact on individual party. The noise introduced to each party's stake has a distinct effect on their odds of winning the slot leader election. Given the uniform variance of injected noise, smaller stakeholders might face evident stake distortion or changes in their winning odds for a single election round. Conversely, for larger stakeholders, the change in their winning odds might not be as prominent. Nevertheless, as the added noises possess zero means, for long-term execution, the expected number of election wins aligns with each party's true stake.

## 6.4 Stake Distortion Protocol Design

**Building blocks**. We build the protocols in a hybrid world where the following auxiliary functionalities are avail-

able: (i) There exists a hiding commitment scheme [51] that given randomness r and message x, an algorithm  $Com_r(x)$  commits x to  $com_x$  with opening r; (ii) We utilize pseudorandom functions [39],  $PRF_k(x)$  with input x and the evaluation key k, for deriving DP noises and rely on an unpredictable random oracle,  $\mathcal{F}_{ro}$ , to produce randomness. An instantiation of  $\mathcal{F}_{ro}$  is the random beacon [11]; (iii) There exists a non-interactive zeroknowledge functionality [39],  $\mathcal{F}_{nizk}^{\mathcal{L}}$ , that allows proving of statements in an NP language  $\mathcal{L}$ ; (iv) There is a "determine stake" functionality,  $\mathcal{F}_{stk}$ , allows parties to access their stake, committed stakes, and stake commitment opening at any time. Additionally,  $\mathcal{F}_{stk}$  allows for querying a list, L, containing: (a) all registered users with their respective committed stake, and (b) each party's noise generation key. In practice, the aforementioned information can be efficiently computed from the blockchain [29].

Timer protocol  $\prod_{\text{timer}}$ . We provide the design of timer protocol  $\prod_{\text{timer}}$  that securely realizes  $\mathcal{F}_{\text{Timer}}$ . Initially, we consider all parties agree on a global clock and each party derives a pair of noise generation keys  $(n_{pk}, n_{sk})$  upon joining the system. Such key pair is sampled by selecting a random private key  $n_{sk}$  and setting  $n_{pk} \leftarrow \text{PRF}_{n_{sk}}(0)$ . The parties manage  $n_{sk}$  themselves, which remains unknown to others, but  $n_{pk}$  is accessible to other parties (i.e., through  $\mathcal{F}_{stk}$ ). Parties can join at any time, while stake distortion occurs every T steps, thus we restrict each new party to wait until the next distortion schedule to acquire the first distorted stake before it can participate in the ledger maintenance. Protocol 6.2 shows details from the viewpoint of party P with identifier pid.

```
Upon receiving DISTORT from pid

1: Obtain current time j from the global clock

2: if j mod T = 0 then (s_j, com_{s_j}, r) \leftarrow \mathcal{F}_{stk}

3: (ra, rb) \leftarrow \mathcal{F}_{ro}(j), (z^0, z^1) \leftarrow PRF_{n_{sk}}(ra, rb)

4: \tilde{s}_j \leftarrow s_j + \frac{\alpha}{\epsilon}(\ln(z^0) - \ln(z^1))

5: Sample \tilde{r} and com_{\tilde{s}_j} \leftarrow Com_{\tilde{r}}(\tilde{s}_j)

6: Get proof: \pi_j \leftarrow \mathcal{F}_{nizk}^{\mathcal{L}_{timer}}(prove, \mathbf{x}, \mathbf{w})

7: return \tilde{s}_j broadcast (pid, com_{\tilde{s}_j}, \pi_j)

8: else \tilde{s}_j = \tilde{s}_{j-1} return \tilde{s}_j

Upon receiving (GET\_COMM, pid)

9: Get (com_{\tilde{s}_j}, \pi_j) for pid from the network.

10: if accept \leftarrow \mathcal{F}_{nizk}^{\mathcal{L}_{timer}}(verify, \pi_j) then

11: Record com_{\tilde{s}_j} for pid
```

In order to distort the stake,  $\prod_{\mathsf{Timer}}$  first acquires necessary information from  $\mathcal{F}_{\mathsf{stk}}$  (or from the blockchain) and obtains the current public randomness from  $\mathcal{F}_{\mathsf{ro}}$ . The protocol then evaluates a PRF to derive private random seeds  $z^2$ ,  $z^0$ ,  $z^1$ , which are subsequently converted into

<sup>&</sup>lt;sup>2</sup>We treat random seeds as fixed-point values in (0, 1).

an instance of Laplace noise drawn from Lap( $\frac{\alpha}{\epsilon}$ ) [50]. Given that n<sub>sk</sub> remains concealed from other parties, no entity can gain knowledge about the derived DP noise. Furthermore, provided the output of  $\mathcal{F}_{ro}$  is unpredictable, the adversary is unable to adaptively select n<sub>sk</sub> to maximize the derived DP noise. Lastly,  $\prod_{Timer}$  returns the distorted stake to party P and broadcasts a hiding commitment to the distorted stake, accompanied by a zeroknowledge proof demonstrating that the committed stake has been accurately distorted to all other parties. Next, we define the statements by their corresponding NP languages: A tuple  $(\mathbf{x}, \mathbf{w}) \in \mathcal{L}_{timer}$  if all the following holds:

- Instance:  $\mathbf{x} = (\epsilon, \alpha, \mathsf{n}_{\mathsf{pk}}, \mathsf{com}_{s_i}, \mathsf{com}_{\tilde{s}_i}, \mathsf{ra}, \mathsf{rb})$
- Witness:  $\mathbf{w} = (s_i, \tilde{s}_i, \mathsf{n}_{\mathsf{sk}}, z^0, z^1, r, \tilde{r})$
- Correct noise generation key:  $PRF_{n_{ek}}(0) = n_{pk}$
- Correct randomness:  $PRF_{n_{sk}}(ra, rb) = (z^0, z^1)$
- Correct stake distortion:  $\tilde{s}_j = s_j + \frac{\alpha}{\epsilon} (\ln(z^0) \ln(z^1))$
- $\exists r, \tilde{r} \text{ s.t. } com_{\tilde{s}_i} = Com_r(\tilde{s}_i) \text{ and } com_{\tilde{s}_i} = Com_{\tilde{r}}(\tilde{s}_i)$

For any other times between two distortion schedules, the protocol simply reuses the distorted stake from the last distortion round. Furthermore, when a party receives a broadcast stake commitment from others, the protocol interacts with  $\mathcal{F}_{\mathsf{nizk}}^{\mathcal{L}}$  to verify the validity of  $\pi$ , accepting the commitment only if  $\pi$  is deemed valid.

**Binary protocol**  $\prod_{Bin}$ . To implement  $\prod_{Bin}$  that securely realizes  $\mathcal{F}_{Bin}$ , we utilize the same strategy as  $\prod_{Timer}$  to obtain randomness and transform DP noises. The tricky part is dealing with the DP interval tree and deriving noisy p-sums. In our design, we use the binary representation of time to implicitly track the tree structure [16], which simplifies the NP statements for generating proofs. The details are provided in Protocol 6.3

Phase one distortion in our design follows the same implementation as  $\prod_{Timer}$ , except the noisy stake is cached in tx\*. In phase two, the protocol derives a new noisy p-sum every T time steps based on the current time's binary representation (6.3:5-11).  $\prod_{Bin}$  then aggregates a subset of previously computed (cached) noisy p-sums to determine the noisy stake (6.3, 12). Note that,  $\prod_{\text{Timer}}$  can safely recycle a subset of "old" p-sums (6.3:10), requiring it to track at most  $\log(\frac{L}{T})$  noisy p-sums. Similarly,  $\prod_{Bin}$  generates proof for validating the stake distortion process, with NP statements defined as follows: A tuple  $(\mathbf{x}, \mathbf{w}) \in \mathcal{L}_{bin}$  if all the following holds:

- $\mathbf{x} = (\epsilon, \alpha, \mathsf{n}_{\mathsf{pk}}, \mathsf{ra}, \mathsf{rb}, \mathsf{com}_{s_i}, \mathsf{com}_{s_{i-T}}, \mathsf{com}_{\tilde{s}_i}, \mathsf{com}_{\mathsf{tx}}, \mathsf{com}_{\tilde{\mathsf{tx}}})$
- $\mathbf{w} = (s_i, s_{i-T}, \tilde{s}_i, \mathsf{tx}, \tilde{\mathsf{tx}}, \mathsf{n}_{\mathsf{sk}}, z^0, z^1, r, \tilde{r}, \overrightarrow{r})$
- $PRF_{n_{sk}}(0) = n_{pk}$  and  $PRF_{n_{sk}}(ra, rb) = (z^0, z^1)$
- Correct new p-sum:  $tx[\ell] = \sum_{k=0}^{\ell} tx[k] + s_j s_{j-T}$
- Noisy p-sum:  $t\tilde{x}[\ell] \leftarrow tx[\ell] + \frac{\alpha}{\epsilon}(\ln(z^0) \ln(z^1))$
- Correct noisy stake:  $\tilde{s}_j \leftarrow \mathsf{tx}^* + \sum_{k: \mathsf{bin}_k(t)=1} \mathsf{tx}[k]$
- $\exists$  opening  $r, \tilde{r}$  s.t. opens  $com_{s_i}, com_{\tilde{s}_i}$  to  $s_i, \tilde{s}_i$ .
- $\forall_k (\exists r_k, \tilde{r}_k \text{ s.t. opens } com_{tx_k}, com_{\tilde{t}x_k} \text{ to } tx[k], tx[k])$

# Protocol 6.3: Binary protocol ∏<sub>Bin</sub>

//vectors tx, and tx are initiated as 0. Upon receiving DISTORT from pid

- 1: Obtain current time *j* from the global clock.
- 2: **if**  $j \mod L = 0$  **then**
- Follow  $\prod_{\text{Timer}}$  steps 2-10, and cache  $\tilde{s}_i$  as  $tx^*$ 3:
- 4: **if**  $t \mod T = 0$ , where  $t \leftarrow j \mod L$  **then**
- $(s_j, \mathsf{com}_{s_i}, r), (s_{j-T}, \mathsf{com}_{s_{j-T}}, r') \leftarrow \mathcal{F}_{\mathsf{stk}}$ 5:
- $(ra, rb) \leftarrow \mathcal{F}_{ro}(j), (z^0, z^1) \leftarrow \mathsf{PRF}_{\mathsf{n}_{sk}}(ra, rb)$ 6:
- Express in binary form:  $t = \sum_{k} 2^{k} \cdot bin_{k}(t)$ 7:
- $\ell \leftarrow \min\{k : \text{bin}_k(t) \neq 0\}$ 8:
- $\mathsf{tx}[\ell] \leftarrow \sum_{k < \ell} \mathsf{tx}[k] + (s_i s_{i-T})$ 9:
- 10: **for**  $k < \ell$  override  $\tilde{\mathsf{tx}}[k] = \mathsf{tx}[k] = 0$
- $\tilde{\mathsf{tx}}[\ell] = \mathsf{tx}[\ell] + \frac{\alpha}{\epsilon} (\mathsf{ln}(z^0) \mathsf{ln}(z^1))$ 11:
- $\tilde{s}_j \leftarrow \mathsf{tx}^* + \sum_{k: \mathsf{bin}_k(t)=1} \tilde{\mathsf{tx}}[k];$ 12:
- Sample  $\tilde{r}$  and  $com_{\tilde{s}_i} \leftarrow Com_{\tilde{r}}(\tilde{s}_i)$ 13:
- Sample  $r_{\ell}$  and  $\mathsf{com}_{\mathsf{tx}_k} \leftarrow \mathsf{Com}_{r_{\ell}}(\mathsf{tx}[\ell])$ 14:
- $\begin{array}{l} \text{Sample } \tilde{r}_{\ell} \text{ and } \mathsf{com}_{\tilde{\mathsf{tx}}_{\ell}} \leftarrow \mathsf{Com}_{\tilde{r}_{\ell}}(\tilde{\mathsf{tx}}[\ell]) \\ \text{Get proof: } \pi_{j} \leftarrow \mathcal{F}_{\mathsf{nizk}}^{\mathcal{L}_{\mathsf{bin}}}(\mathsf{prove}, \mathbf{x}, \mathbf{w}) \end{array}$ 15:
- 16:
- 17:
- **broadcast** (pid,  $com_{\tilde{s}_i}$ ,  $com_{tx_\ell}$ ,  $com_{\tilde{tx}_\ell}$ ,  $\pi_i$ ) 18:
- 19: **else**  $\tilde{s}_i = \tilde{s}_{i-1}$  **return**  $\tilde{s}_i$

Upon receiving (GET\_COMM, pid)

- 20: Get  $com_{\tilde{s}_i}$ ,  $com_{tx_\ell}$ ,  $com_{\tilde{tx}_\ell}$ ,  $\pi_i$  for pid
- 21: Override  $\forall_{k<\ell} \operatorname{com}_{\mathsf{tx}_k} = \operatorname{com}_{\mathsf{tx}_k} = \operatorname{Com}(0)$  for pid
- 22: **if** accept  $\leftarrow \mathcal{F}_{\text{nizk}}^{\mathcal{L}_{\text{bin}}}(\text{verify}, \pi_j)$  **then**
- 23: Record  $com_{\tilde{s}_j}$ ,  $com_{tx_{0:\ell}}$ , and  $com_{t\tilde{x}_{0:\ell}}$  for pid

Observe that  $\prod_{Bin}$  requires verifying nodes to cache additional objects per party, i.e., the commitments to  $tx, \tilde{tx}$ . However, one needs to keep track of at most  $\log(\frac{L}{T})$  psums(commitments). Hence, the storage blowup in contrast to  $\prod_{\mathsf{Timer}}$  is bounded by  $O(\log(\frac{L}{T}))$ .

**Theorem 9.**  $\prod_{\mathsf{Timer}} and \prod_{\mathsf{Bin}} UC \ emulates \ [15] \mathcal{F}_{\mathsf{Timer}}$ and  $\mathcal{F}_{Bin}$ , respectively, in the  $(\mathcal{F}_{stk}, \mathcal{F}_{ro}, \mathcal{F}_{nizk}^{\perp})$ -hybrid world with the presence of a p.p.t. adversary.

*Proof.* Please refer to the Appendix A

# Case Study with Real-World System

In this section, we perform case studies on our stake distortion protocols using a real-world PoS blockchain, Ethereum 2.0 [27], and examine its impact on original guarantees. We explore the practicality of stake distortion by addressing several key questions.

- Q1. Will the ledger still provide safety guarantees after adopting stake distortion? What privacy level can be achieved given certain safety requirements?
- Q2. How does  $\mathcal{F}_{Timer}$  compare to  $\mathcal{F}_{Bin}$  in privacy decay? How one can benefit from a smoother decay?
- Q3. Does stake distortion mitigate the SIA risks?

Total stake, S	13,488,174 Eth	
Minimum stake, $v$	32 Eth	
Slot time (delay)	12 seconds	
Distortion interval, T	4 days	
$\mathcal{F}_{Bin}$ Phase 1 interval, $L$	$45 \times T$	

Table 1: Ethereum 2.0 parameters setup

**Setup.** We implement simulations based on Ethereum 2.0's phase 0 specification [27] and utilize key statistics from the beacon chain [28]. With an average daily transaction volume per unique address below 0.25, it's reasonable to assume that Ethereum 2.0 parties' stakes remain stable over a 4-day period. Consequently, we set the distortion schedule T = 4 days and  $\mathcal{F}_{Bin}$ 's phase 1 interval L = 45T (6 months). Table 1 summarizes the setup.

# 7.1 Safety Guarantee

We plot safety curves in Figure 4 based on Theorem 8 with  $\delta = 10^{-9}$  (safety failure probability) 3 and assume  $\tilde{f} < \frac{1}{3}$  as the general safety requirement. Given  $\alpha$  (resp.  $\epsilon$ ), the additive bound for distinguishing two stake values, the curves imply the minimum privacy loss,  $\epsilon$  (resp. the maximum  $\alpha$ ), achievable at certain safety levels. For each mechanism, we plot three curves, which correspond 20%, 25%, and 30% of malicious tolerance, respectively.

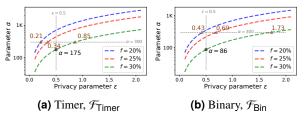
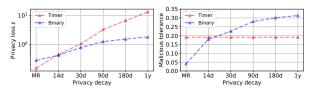


Figure 3: Safety upper bounds

Observation-1. Both mechanisms provide reasonable stake privacy while retaining safety at the same time. In general, even with a relatively small  $\epsilon$  and a large malicious threshold, i.e.,  $\epsilon = 0.5$  [4] and f = 30%,  $\alpha$  can be set up to 175 and 86, respectively, for  $\mathcal{F}_{\text{Timer}}$  and  $\mathcal{F}_{\text{Bin}}$ . By the latest blockchain data [20], only less than 0.05% of accounts have more than 80 Eth, which suggests that privacy is ensured for at least 99.95% of total users.

Observation-2. There is a trade-off between privacy and safety guarantees. Apparently, one can achieve a larger  $\alpha$  or a smaller privacy loss when considering a relatively weaker safety (smaller malicious tolerance). This suggests a trade-off between privacy and safety which provide flexibility for practitioners to better configure the protocols. For example, if there is a slashing mechanism [2] to help reduce the malicious fraction, one may choose to gain better privacy by trading off some safety.



(a) Privacy loss vs. Time

(b) Safety vs. privacy decay

Figure 4: Safety upper bounds

Given that the two mechanisms have different privacy decay, to better study them, we compare  $\mathcal{F}_{Timer}$  and  $\mathcal{F}_{Bin}$  under different privacy decay goals, i.e., privacy loss for the most recent (MR) time, the last 14, 30, 90, 180 and 365 days. First, we bound the malicious tolerance, i.e. 20%, and set  $\alpha = 200$  then plot their privacy loss over time in Figure 4a. Next, we compare them in a different angle in Figure 4b, where we keep the configurations of  $\mathcal{F}_{Timer}$  unchanged but modify  $\mathcal{F}_{Bin}$  to enforce it reaches the same privacy loss as  $\mathcal{F}_{Timer}$  under the different privacy decay goals. We then compare the malicious tolerance (safety) ensured by two mechanisms.

Observation-3. For short-term privacy,  $\mathcal{F}_{Timer}$  outperforms  $\mathcal{F}_{Bin}$ , while the result is reversed when focusing on long-term privacy. As shown in Figure [4a]  $\mathcal{F}_{Timer}$  exhibits lower short-term privacy loss than  $\mathcal{F}_{Bin}$  when fixing the safety level. This occurs since  $\mathcal{F}_{Timer}$  adds a single instance of DP noise, leading to smaller adversary stake growth and enabling the configuration under a smaller  $\epsilon$ . However, due to logarithmic privacy decay,  $\mathcal{F}_{Bin}$  offers better long-term privacy loss even with a larger  $\epsilon$ . Figure [4b] shows that if one is more concerned about long-term privacy, it is better to employ  $\mathcal{F}_{Bin}$  as it ensures a larger malicious tolerance, which up to  $1.63 \times$  that of  $\mathcal{F}_{Timer}$ . Overall, this observation provides practitioners with a sense of how they would benefit from a stake distortion that has a smoother privacy decay.

# 7.2 Safety Simulation

To further validate the safety bounds we conduct stake distortion simulations. Specifically, we simulate a set of stakeholders and assume the attacker corrupts different fractions of them. Subsequently, we configure  $\epsilon$ ,  $\alpha$  as per Theorem 8 and run stake distortion over each party. Finally, we investigate whether the corrupted stake after distortion,  $\tilde{f}$ , exceeds 1/3. Additionally, to max out the attacker's stake increment, we consider each individual stakeholder to hold a minimum staking, and thus the number of parties is maxed out as  $n = \frac{S}{v}$ . Table 2 summarizes the result over 10K repeated runs.

**Observation-4.** With proper setups as per the safety bounds,  $\hat{f} < \frac{1}{3}$  with overwhelming probability. According to Table 2 the corrupted stake,  $\tilde{f}$ , after distortion for all testing groups does not exceed 0.3314, which

<sup>&</sup>lt;sup>3</sup>Similar to a practical bound chosen by Algorand [17]

 $<sup>^4</sup>$ In comparison with Google community mobility report which uses  $\epsilon=2.4$  per user per day,  $\epsilon=0.5$  is quite small.

	Timer $\mathcal{F}_{Timer}$		Binary $\mathcal{F}_{Bin}$	
f (%)	$(\epsilon, \alpha)$	$\max  ilde{f}$	$(\epsilon, \alpha)$	$\max  ilde{f}$
10%	(0.5, 1214)	0.33082	(0.5, 552)	0.33078
15%	(0.5, 963)	0.32952	(0.5, 475)	0.32869
20%	(0.5, 701)	0.33055	(0.5, 346)	0.33132
25%	(0.5, 438)	0.32858	(0.5, 216)	0.32336
30%	(0.5, 175)	0.32309	(0.5, 86)	0.32626

Table 2: Simulation of stake distortion

is consistent with the primary safety requirement, i.e.,  $\tilde{f} < 1/3$ . Consider that stake distortions occur every 4 days, thus the 10K repeated runs can actually simulate a 109-year operation of the blockchain coupled with stake distortion. The simulation implies in those 109 years of operation, there has not been even one violation of the safety requirement. In fact, as  $\beta$  is set to  $10^{-9}$ , the probability of even one violation happens should smaller than  $10^{-4}$  (taking union bound). In general, we conclude that the simulations validate our safety upper bounds.

#### 7.3 Resilience to SIA

We investigate whether stake distortion helps mitigate SIA risk, which is achieved by initiating RdBin on simulated Ethereum 2.0's ledger maintenance protocol with and without stake distortion, then compare their relative inference errors. Moreover, we are also interested in the SIA resistance of  $\mathcal{F}_{\text{Timer}}$  and  $\mathcal{F}_{\text{Bin}}$  under different safety requirements. To address these, we simulate the protocol for  $1.29 \times 10^6$  steps (6 months in the real world) and issue one RdBin attack every 4 days. By default, we configure RdBin with  $\theta = v$ ,  $\tau = 0.01$ , and run the attack against a target with random stake  $f_v \in (0, f)$ . When simulating Ethereum 2.0 with stake distortion, we assume the protocol tolerates different malicious fractions, i.e., from 10% to 30%, configure related parameters accordingly, and run attacks against each group independently. We report the relative inference errors of RdBin in Figure 5.

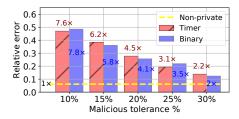


Figure 5: Privacy loss vs. Time

Observation-5. Both  $\mathcal{F}_{Timer}$  and  $\mathcal{F}_{Bin}$  mitigate SIA risk, and the simulation result suggests a trade-off between safety guarantees and SIA resistance. According to Figure 5, when stake distortion is employed, the relative SIA inference error is at least twice that of the standard protocol and can reach up to 7.8×. This further implies the practicability of stake distortion in mitigating SIA risks. Moreover, the relative error of both  $\mathcal{F}_{Timer}$  and  $\mathcal{F}_{Bin}$ decreases when the malicious tolerance increases, i.e., the

		Timer	Binary
Setup	Setup time	19.64 s	20.63 s
	Proving key size	38 Mb	38 Mb
	Verification key size	4 Kb	4 Kb
Prove	Proof time	7.07 s	7.32 s
	Proof size	4 Kb	4 Kb
Verify	Verification time	11 ms	10ms
Other	Time	1.33 s	1.89 s
	Storage	22.72 Mb	62.21 Mb

Table 3: Evaluation result for the prototype protocol weaker the safety demanded, the stronger the resistance to SIA obtained, and vice versa.

#### 7.4 **Key Takeaways**

The first two observations address Q1, where we learn that our mechanisms provide reasonable privacy while preserving safety (liveness). Moreover, we are aware of a privacy-safety tradeoff in our mechanisms. Observation 3 compares the two mechanisms and addresses Q2, from which we learn that one should choose  $\mathcal{F}_{Timer}$  when focusing on short-term privacy goals, while  $\mathcal{F}_{Bin}$  is a better option for long-term goals. The safety simulations (Section 7.2) validate our theoretical bounds while the SIA experiments (Section 7.3) address Q3 and demonstrate our mechanisms indeed mitigate SIA risks.

## **Performance Evaluation**

Implementation and configuration. We utilize the same method as [51] to implement PRF and Com using SHA256 compression functions. Fractional numbers are stored as 32-bit fixed-point values with a scaling factor of 2<sup>16</sup>, and we consider the public randomness as 16-bit unsigned fixed-point values with a scaling factor of  $2^{16}$ , ensuring the randomness is spread across (0,1). We use the Remez algorithm [53] to create approximated polynomials for the ln(x) circuit. Alternative techniques [48,52] and engineering optimizations exist for implementing the ln(x) circuit, potentially improving approximation accuracy or computational efficiency. Developing optimized protocols is not the primary focus of this work, but may be of independent interest. Nevertheless, even without these optimizations, our evaluations still show a reasonable overhead associated with the stake distortion protocols. To implement the NIZK proof, we utilize zk-SNARK library Zokrates-0.8.5 (Rust implementation) [1] under the proving scheme Groth16 [34]. In addition, we assume the existence of the parties' stake commitments and their integrity have been verified, as in practice such commitments can be effectively obtained and verified [51] from the blockchain. The protocols are evaluated over machines with 2.6GHz CPU and 16Gb RAM.

**Evaluation result.** Table 3 summarizes the evaluation results, apparently, the performance overhead is dominated by the costs associated with the NIZK proof. Nevertheless, the adoption of our protocol may not impose an excessive performance overhead, for example, Zcash requires more than 120 seconds of proof time (without optimization) per block [51]. Thus, adopting our protocols only increase the proof overhead by at most 5.8%. Furthermore, since stake distortion doesn't transpire at every time step, the amortized overhead may be reduced. For example, assuming stake distortion occurs every 4 days, the amortized overhead in the Zcash case could be less than 0.002%. Note that stake distortion proofs can be combined with other ledger-specified proofs, enabling the verifier to validate all SNARK arguments via a single verification. For the storage overhead, both protocols require parties to cache additional data objects, i.e., noisy stake, commitments, etc. However, the storage volume does not exceed 63 Mb, which represents a slight overhead compared to storing the full blockchain data (i.e., 895 Gb for Ethereum [28]).

# **Discussion**

Mitigating privacy decay. One can create a stake distortion strategy without cumulative privacy loss over time. For example, for every T time steps, the party generates a noisy p-sum representing all transaction outcomes from the last T steps, and the noisy stake is derived by summing up all released noisy p-sums. This approach guarantees that distortion takes place across distinct transactions, and the extraction of the noisy stake is simply a postprocessing step involving noisy p-sums. As per parallel composition and post-processing theorems [26] of DP, the total privacy loss is bounded by  $\epsilon$ . Although privacy loss doesn't accumulate, injected noise variance does, leading to reduced safety guarantees over time, i.e., the upper bound of adversarial stake increment increases. Another alternative is adopting the relaxed  $(\epsilon, \delta)$ -DP guarantee, allowing unbounded privacy loss with a small failure probability  $\delta$ . This allows injecting Gaussian noises (offer tighter composition bounds than Laplace noises) that optimize privacy decay.

Permissionless clock synchronization. This work considers the globally synchronized setting in line with SOTA designs [40, 19, 39] such that all parties can retrieve time from a global clock (functionality). We say that one can also adopt permissionless clock synchronization protocols [4], i.e. if some (honest) party believes the global time as j, all parties believe it as  $j \pm \delta_c$ , with a small  $\delta_c$ , if a global clock is not accessible. Although, there might be small chances for a party that fails the slot leader election due to the use of outdated noisy stake, i.e. within  $\delta_c$  right after a (globally defined) distortion schedule. Since  $\delta_c$  is typically small [4] and the stake distortion interval T is relatively large, thus for the majority of time slots, all parties agree on the most recent distorted stake. While integrating our design with permissionless clock synchronization is not this paper's focus, it is considered an important future enhancement.

#### **Related Works** 10

Proof-of-Stake private ledgers. The first PoS design appears in [42]. Followed by this, there has been a series of efforts [9, 17, 10, 40, 19, 3, 31, 27, 58, 14, 30, 60] on formulating PoS models as well as providing protocols with provable security guarantees. Recently, several works have been proposed to address the stake and transaction privacy of PoS ledgers, which include Ouroborous Crypsinous [39], Ganesh et al. [29], and Baldimtsi et al. [6]. However, the stake privacy is typically ensured at the cost of assuming complete anonymity of parties' identity [29, 39]. Unfortunately, [43] demonstrates the existence of a tension between liveness and anonymity. In this work, we provide the first-of-its-kind solution that provides provable stake privacy for PoS private ledgers in a manner that does not assume anonymity.

Stake inference attacks. Both [39] and [29] state that the parties' stake information may be revealed with the execution of the PoS protocols, no specific attacks for inferring stake have been proposed until the publication of [43]. However, their proposed attack considers only deterministic protocols and requires linear time complexity. In this work, we provide practical stake inference attack that is proven to be valid against randomized protocols with probabilistic liveness. Moreover, our attack only yields sublinear time complexity.

Differential privacy. Differential privacy (DP) introduced by Dwork et.al. [26] is currently the "de-facto" standard for achieving data privacy. Since 2010, couple of efforts have been proposed to address DP under dynamic setting [25, 24, 16, 38, 55, 56]. Among these works, [12, 38] employs a relaxed privacy with decay model to ensure better utility, which is similar to our formulation. Nevertheless, as far as we know, we are the first to incorporate DP with private PoS blockchains.

#### 11 Conclusion

In this paper, we present the first practical stake inference attack against randomized PoS protocols. The existence of such an attack further implies the inadequacy of SOTA PoS designs in achieving stake and transaction privacy. To formulate a rigorous privacy definition, we incorporate DP with the standard UC definition of private ledgers. Guided by the derived privacy model, we design DP stake distortion protocols that assist existing PoS protocols in resolving stake and transaction privacy.

## Acknowledgments

This paper is supported in part by NSF Awards 2016393 and 2237814.

## References

- [1] A toolbox for zksnarks on ethereum. https://github.com/ Zokrates/ZoKrates, 2022.
- [2] Aumasson, J.-P., Kolegov, D., and Stathopoulou, E. Security review of ethereum beacon clients. arXiv preprint arXiv:2109.11677 (2021).
- [3] BADERTSCHER, C., GAŽI, P., KIAYIAS, A., RUSSELL, A., AND ZIKAS, V. Ouroboros genesis: Composable proof-of-stake blockchains with dynamic availability. In Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (2018), pp. 913–930.
- [4] BADERTSCHER, C., GAZI, P., KIAYIAS, A., RUSSELL, A., AND ZIKAS, V. Ouroboros chronos: Permissionless clock synchronization via proof-of-stake. Cryptology ePrint Archive (2019).
- [5] BADERTSCHER, C., MAURER, U., TSCHUDI, D., AND ZIKAS, V. Bitcoin as a transaction ledger: A composable treatment. In Annual international cryptology conference (2017), Springer, pp. 324-
- $[6] \quad Baldimtsi, F., Madathil, V., Scafuro, A., and Zhou, L. Anony$ mous lottery in the proof-of-stake setting. In 2020 IEEE 33rd Computer Security Foundations Symposium (CSF) (2020), IEEE, pp. 318-333.
- [7] Bellare, M. Lectures on nizks: A concrete security treatment.
- [8] Bentov, I., Gabizon, A., and Mizrahi, A. Cryptocurrencies without proof of work. In International conference on financial cryptography and data security (2016), Springer, pp. 142–157.
- [9] Bentov, I., Lee, C., Mizrahi, A., and Rosenfeld, M. Proof of activity: Extending bitcoin's proof of work via proof of stake [extended abstract] y. ACM SIGMETRICS Performance Evaluation Review 42, 3 (2014), 34-37.
- [10] Bentov, I., Pass, R., and Shi, E. Snow white: Provably secure proofs of stake. IACR Cryptol. ePrint Arch. 2016, 919 (2016).
- [11] Bhat, A., Kate, A., Nayak, K., and Shrestha, N. Optrand: Optimistically responsive distributed random beacons. Cryptology ePrint Archive (2022).
- [12] Bolot, J., Fawaz, N., Muthukrishnan, S., Nikolov, A., and TAFT, N. Private decayed predicate sums on streams. In Proceedings of the 16th International Conference on Database Theory (2013), pp. 284-295.
- [13] Bowe, S., Chiesa, A., Green, M., Miers, I., Mishra, P., and Wu, H. Z: Enabling decentralized private computation. In 2020 IEEE Symposium on Security and Privacy (2019), pp. 947–964.
- [14] BUCHMAN, E. Tendermint: Byzantine fault tolerance in the age of blockchains. PhD thesis, University of Guelph, 2016.
- [15] CANETTI, R. Universally composable security: A new paradigm for cryptographic protocols. In Proceedings 42nd IEEE Symposium on Foundations of Computer Science (2001), IEEE.
- [16] CHAN, H., SHI, E., SONG, D., ET AL. Private and continual release of statistics. In International Colloquium on Automata, Languages, and Programming (2010), Springer, pp. 405-417.
- [17] CHEN, J., AND MICALI, S. Algorand. arXiv preprint arXiv:1607.01341 (2016).
- [18] Croman, K., Decker, C., Eyal, I., Gencer, A. E., Juels, A., Kosba, A., Miller, A., Saxena, P., Shi, E., Gün Sirer, E., ET AL. On scaling decentralized blockchains. In International conference on financial cryptography and data security (2016), Springer, pp. 106-125.
- [19] DAVID, B., GAŽI, P., KIAYIAS, A., AND RUSSELL, A. Ouroboros praos: An adaptively-secure, semi-synchronous proof-of-stake blockchain. In Annual Conference on the Theory and Applications of Cryptographic Techniques (2018), Springer, pp. 66-98.

- [20] DAY, A., MEDVEDEV, E., NIRMAL, A., AND PRICE, W. Introducing six new cryptocurrencies in bigquery public datasets-and how to analyze them. Google Cloud (2021).
- [21] Delgado-Segura, S., Bakshi, S., Pérez-Solà, C., Litton, J., PACHULSKI, A., MILLER, A., AND BHATTACHARJEE, B. Txprobe: Discovering bitcoin's network topology using orphan transactions. In International Conference on Financial Cryptography and Data Security (2019), Springer, pp. 550-566.
- [22] Dodge, Y., and Cox, D. The Oxford dictionary of statistical terms. Oxford University Press, USA, 2003.
- [23] DOUCEUR, J. R. The sybil attack. In International workshop on peer-to-peer systems (2002), Springer, pp. 251-260.
- [24] DWORK, C. Differential privacy in new settings. In Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms (2010), SIAM, pp. 174-183.
- [25] DWORK, C., NAOR, M., PITASSI, T., AND ROTHBLUM, G. N. Differential privacy under continual observation. In Proceedings of the forty-second ACM symposium on Theory of computing (2010), pp. 715-724.
- [26] DWORK, C., ROTH, A., ET AL. The algorithmic foundations of differential privacy. Found. Trends Theor. Comput. Sci. 9, 3-4 (2014), 211-407.
- [27] ETHERUM. Etherum 2.0 phse 0 becaon chain. //github.com/ethereum/consensus-specs/blob/dev/ specs/phase0/beacon-chain.md 2021.
- [28] ETHSCAN. Ethscan Statistics. https://ethscan.org/ 2022.
- [29] GANESH, C., ORLANDI, C., AND TSCHUDI, D. Proof-of-stake protocols for privacy-aware blockchains. In Annual International Conference on the Theory and Applications of Cryptographic Techniques (2019), Springer, pp. 690-719.
- [30] GAO, Y., AND NOBUHARA, H. A proof of stake sharding protocol for scalable blockchains. Proceedings of the Asia-Pacific Advanced Network 44, 1 (2017), 13-16.
- [31] Gaži, P., Kiayias, A., and Zindros, D. Proof-of-stake sidechains. In 2019 IEEE Symposium on Security and Privacy (SP) (2019), IEEE, pp. 139-156.
- [32] Grinstead, C. M., and Snell, J. L. Introduction to probability. American Mathematical Soc., 1997.
- [33] Groth, J. Simulation-sound nizk proofs for a practical language and constant size group signatures. In International Conference on the Theory and Application of Cryptology and Information Security (2006), Springer, pp. 444-459.
- [34] Groth, J. On the size of pairing-based non-interactive arguments. In Advances in Cryptology-EUROCRYPT 2016: 35th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Vienna, Austria, May 8-12, 2016, Proceedings, Part II 35 (2016), Springer, pp. 305–326.
- [35] Haney, S., Machanavajjhala, A., Abowd, J. M., Graham, M., Kutzbach, M., and Vilhuber, L. Utility cost of formal privacy for releasing national employer-employee statistics. In Proceedings of the 2017 ACM International Conference on Management of Data (2017), pp. 1339-1354.
- [36] HELLMAN, M., AND RAVIV, J. Probability of error, equivocation, and the chernoff bound. IEEE Transactions on Information Theory 16, 4 (1970), 368-372.
- [37] KARP, R. M., AND KLEINBERG, R. Noisy binary search and its applications. In Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms (2007), pp. 881-890.
- [38] Kellaris, G., Papadopoulos, S., Xiao, X., and Papadias, D. Differentially private event sequences over infinite streams. Proceedings of the VLDB Endowment 7, 12 (2014), 1155-1166.

- [39] KERBER, T., KIAYIAS, A., KOHLWEISS, M., AND ZIKAS, V. Ouroboros crypsinous: Privacy-preserving proof-of-stake. In 2019 IEEE Symposium on Security and Privacy (SP) (2019), IEEE, pp. 157–174.
- [40] KIAYIAS, A., RUSSELL, A., DAVID, B., AND OLIYNYKOV, R. Ouroboros: A provably secure proof-of-stake blockchain protocol. In *Annual international cryptology conference* (2017), Springer, pp. 357–388.
- [41] KIFER, D., AND MACHANAVAJJHALA, A. Pufferfish: A framework for mathematical privacy definitions. ACM Transactions on Database Systems (TODS) 39, 1 (2014), 1–36.
- [42] KING, S., AND NADAL, S. Ppcoin: Peer-to-peer crypto-currency with proof-of-stake. self-published paper, August 19, 1 (2012).
- [43] KOHLWEISS, M., MADATHIL, V., NAYAK, K., AND SCAFURO, A. On the anonymity guarantees of anonymous proof-of-stake protocols. In 2021 IEEE Symposium on Security and Privacy (SP) (2021), IEEE, pp. 1818–1833.
- [44] KRAFFT, O., AND SCHMITZ, N. A note on hoeffding's inequality. Journal of the American Statistical Association 64, 327 (1969), 907–912.
- [45] MAZIERES, D. The stellar consensus: A federated model for internet-level consensus. Stellar Development Foundation (2015).
- [46] MILLER, A., LITTON, J., PACHULSKI, A., GUPTA, N., LEVIN, D., SPRING, N., AND BHATTACHARJEE, B. Discovering bitcoin's public topology and influential nodes. *et al* (2015).
- [47] NAKAMOTO, S., AND BITCOIN, A. A peer-to-peer electronic cash system. Bitcoin.—URL: https://bitcoin. org/bitcoin. pdf 4 (2008).
- [48] NARAYAN, A., FELDMAN, A., PAPADIMITRIOU, A., AND HAE-BERLEN, A. Verifiable differential privacy. In *Proceedings of the European Conference on Computer Systems* (2015), pp. 1–14.
- [49] NOETHER, S., MACKENZIE, A., ET AL. Ring confidential transactions. *Ledger 1* (2016), 1–18.
- [50] Ross, S. M. Intro. to probability models. Academic press, 2014.
- [51] SASSON, E. B., CHIESA, A., GARMAN, C., GREEN, M., MIERS, I., TROMER, E., AND VIRZA, M. Zerocash: Decentralized anonymous payments from bitcoin. In 2014 IEEE symposium on security and privacy (2014), IEEE, pp. 459–474.
- [52] SHIRALI, S. A. The bhaskara-aryabhata approximation to the sine function. *Mathematics Magazine* 84, 2 (2011), 98–107.
- [53] TAWFIK, S. A. Minimax approximation and remez algorithm.
- [54] VUJICIC, D., JAGODIC, D., AND RANDIC, S. Blockchain technology, bitcoin, and ethereum: A brief overview. In 2018 17th international symposium infoteh-jahorina (2018), IEEE, pp. 1–6.
- [55] WANG, C., BATER, J., NAYAK, K., AND MACHANAVAJJHALA, A. Dp-sync: Hiding update patterns in secure outsourced databases with differential privacy. In *Proceedings of the 2021 International Conference on Management of Data* (2021), pp. 1892–1905.
- [56] WANG, C., BATER, J., NAYAK, K., AND MACHANAVAJJHALA, A. Incshrink: Architecting efficient outsourced databases using incremental mpc and differential privacy. arXiv preprint arXiv:2203.05084 (2022).
- [57] WANG, C., PUJO, D., NAYAK, K., AND MACHANAVAJJHALA, A. Private proof-of-stake blockchains using differentially-private stake distortion. *Cryptology ePrint Archive* (2023).
- [58] Wood, G. Polkadot: Vision for a heterogeneous multi-chain framework. *White Paper 21* (2016), 2327–4662.
- [59] Xu, B., Luthra, D., Cole, Z., and Blakely, N. Eos: An architectural, performance, and economic analysis. *Retrieved June* 11 (2018), 2019.
- [60] ZAMFIR, V. Casper the friendly ghost: A correct by construction blockchain consensus protocol.

# **A Security Proof**

In this section, we present the primary security proof (for Theorem [9]) related to the stake distortion protocol proposed in this work. To minimize redundancy, we focus on the security proof for the Timer protocol. However, we emphasize that using the same proof technique, one can easily derive a proof for the Binary protocol. Interested readers may refer to our full version [57] for more details.

First, we provide the formal functionality descriptions (1.1) for the timer ( $\mathcal{F}_{Timer}$ ) mechanism. In general, parties can interact with  $\mathcal{F}_{Timer}$  with DISTORT to retrieve the distorted stake at each time, or with GET\_COMM to obtain distorted stake commitments for every other party.

```
Functionality 1.1: \mathcal{F}_{Timer}
        For every time step j evaluate the following:
  1: if j \mod T = 0 then
               for party i = 1, 2, ... do
 2:
                        Determine the true stake s_i^i owned by P_i
 3:
                       \tilde{s}_{i}^{i} \leftarrow s_{i}^{i} + \operatorname{Lap}(\frac{\alpha}{\epsilon})
  4:
                       Commit \tilde{s}_{j}^{i} to com_{\tilde{s}_{j}^{i}} with r_{j}^{i} Record \tilde{s}_{j}^{i}, com_{\tilde{s}_{i}^{i}}, and r_{j}^{i} for P_{i}
  5:
 6:
  7: else
 8:
                       \begin{split} \tilde{s}^i_j \leftarrow \tilde{s}^i_{j-1}, r^i_j \leftarrow r^i_{j-1}, \operatorname{com}_{\tilde{s}^i_j} \leftarrow \operatorname{com}_{\tilde{s}^i_{j-1}} \\ \operatorname{Record} \tilde{s}^i_j, \operatorname{com}_{\tilde{s}^i_i}, \operatorname{and} r^i_j \text{ for } P_i \end{split}
 9:
10:
        Upon receiving (DISTORT, sid) from P_i:
11: return \tilde{s}_{j}^{i}, com_{\tilde{s}_{i}^{i}}, and r_{j}^{i}.
12: add com_{\tilde{s}_{i}^{i}} to every registered party's local cache.
        Upon receiving (GET_COMM, id, sid) from P_i:
13: return com_{\tilde{s}_{i}^{k}} from P_{i}'s cache
```

*Proof.* (**Theorem 9**) We prove this theorem by constructing a simulator  $S_{timer}$ , which corrupts the same set of nodes as the real-world adversary  $\mathcal{A}$  and interacts with  $\mathcal{F}_{sd}$ . This simulator can produce a transcript that is computationally indistinguishable from the real-world protocol execution of  $\prod_{sd}$  in the presence of  $\mathcal{A}$ . We assume the existence of a simulator that simulates the message broadcast channel (the network) as described in [29]. Consequently, we will not delve into the details of simulating message broadcasting in real-world execution. For more information, readers may refer to [29].

To construct the simulator, we consider the existence of a simulator,  $\mathcal{S}_{\text{nizk}}^{\mathcal{L}_{\text{limer}}}$ , for non-interactive zero-knowledge (NIZK) that proves statements in the NP language of  $\mathcal{L}_{\text{timer}}$ .  $\mathcal{S}_{\text{nizk}}^{\mathcal{L}_{\text{limer}}}$  is able to simulate a proof using a trapdoor (without requiring the witness) that is indistinguishable from one an honest prover would provide with the witness [7] [33]. We refer to the proof generated by  $\mathcal{S}_{\text{nizk}}^{\mathcal{L}_{\text{timer}}}$  using the trapdoor as the simulated proof.

Let HYB<sub>0</sub> represent the distribution of real-world protocol execution of  $\prod_{timer}$ . We consider the world HYB<sub>1</sub>, which is the same as the protocol execution, except for the following: the noise generation keys and the DP noises are obtained from the simulator  $S_{timer}$ . Specifically,  $S_{timer}$ samples a random string, using it as the key to deriving the private randomness. By the security of PRF the output of (private randomness)  $PRF_{n_{sk}}(0)$  is computationally indistinguishable from  $PRF_{rd\_string}(0)$ . Furthermore, by the fundamental laws of probability transformation, Laplace or Gaussian random variables that are derived from the actual and the simulated randomness are also indistinguishable. Hence, we say that the distributions of HYB<sub>0</sub> and HYB<sub>1</sub> are indistinguishable.

Next, we consider another world HYB2, which is identical to the world HYB<sub>1</sub> except for replacing the calls to  $\mathcal{F}_{\text{stk}}$  with the interaction to  $\mathcal{S}_{\text{timer}}$ . The simulator  $S_{\text{timer}}$  simulates the honest party's stake commitment using Comr(0), where r is a random opening, and returns the honest party's noise generation key with the simulated one (as demonstrated in  $HYB_1$ , the simulated  $n_{pk}$ is indistinguishable from the actual key). Owing to the equivocality property of commitment schemes [39], we can conclude that the distributions of HYB<sub>1</sub> are indistinguishable from those of HYB<sub>2</sub>.

Lastly, we consider the ideal world HYB3. The only difference between HYB3 and HYB2 is that in HYB3, the noisy stake commitments for honest parties are simulated using (i) 0 stake as the true value; (ii) simulated DP noise z; (iii) a random opening to derive the commitment; and (iv) a simulated NIZK proof generated by  $\mathcal{S}^{\mathcal{L}timer}$ nizk. Due to the equivocality property of commitments, we know that the simulated noisy stake commitments for honest parties are indistinguishable from the actual ones. Furthermore, the simulation security of NIZKs (or zeroknowledge property) ensures that the simulated proof by  $S^{\text{Ltimer}}$ nizk, generated without using a witness, is computationally indistinguishable from the actual proof obtained from  $\mathcal{F}_{nizk}^{\mathcal{L}}$ . Consequently, we conclude that the distribution of the ideal world HYB3 is indistinguishable from that of HYB<sub>2</sub> and, by extension, also indistinguishable from the real-world executions, HYB<sub>0</sub>.

# **Proof of Theorems**

#### **Proof of Theorem 1 B.1**

**Theorem 10.** Given  $\delta \in (0,1)$ , with probability at least  $1-\delta$ ,  $pSC(...,\delta)$  outputs correct comparisons.

*Proof.* We first prove that for each round, the probability that Algorithm 1 outputs the wrong comparison between the two biased coins  $q(f_v)$  and  $q(f_{cmp})$  is at most  $\delta_i$ . Without loss of generality, we consider the case where the algorithm outputs  $f_v > f_{cmp}$   $(p_0 > p_1 \text{ and } |p_0 - p_1| > 2\tau_i)$ . By Hoeffding's Inequality [44] the failure probability that Algorithm 1 outputs the wrong comparison satisfies

Pr [wrong\_cmp]

$$\leq \Pr\left[q(f_{v}) < p_{0} - \tau_{i}\right] + \Pr\left[q(f_{\mathsf{cmp}}) > p_{1} + \tau_{i}\right]$$

$$\leq e^{-2} e^{\log\left(\frac{1}{\delta_{i}}\right)} = e^{-2} \delta_{i}$$
(4)

By union bound, the overall failure probability is bounded by  $\sum_{i=0}^{+\infty} \delta_i < \frac{\delta}{\alpha-1} < \delta$ 

**Theorem 11.** The expected running time of pSC is bounded by  $O\left(\frac{\log(1/\delta) + \log(1/\max(\tau, \tau_c))}{\max(\tau, \tau_c)^2}\right)$  (sim) coin flips.

*Proof.* We prove the complexity of pSC by following the same technique used by [37] for proving Lemma 3.2. Without loss of generality, we consider  $q(f_v) > q(f_{cmp})$ . We start with the case of  $\tau_c > \tau$ , and let  $k = \log(\frac{1}{\tau_c})$ , so for any round  $\ell > k$  we have  $\tau_{\ell} < \frac{\tau_c}{2}$  . By Hoeffding's Inequality 44 the probability such that pSC keeps running after  $\ell > k$  rounds is at most  $\delta \cdot e^{-\left(\frac{1}{2}e^{2\ell}\tau_c^2 + \ell\right)}$ . The running time of round  $\ell > k$  grows exponentially in  $\ell$ , while the probability that Algorithm 11 to continue running after round  $\ell > k$  decreases faster than exponential in  $\ell$ . Hence, the expected running time is dominated by the running time of round k, which is  $O\left(\frac{\log(1/\delta) + \log(1/\tau_c)}{\tau_c^2}\right)$ . Next, we consider the case where  $\tau_c < \tau$ . With the same technique, one can obtain that the expected running time is bounded by the running time of  $k^{th}$  round where  $k = \log(\frac{1}{\tau})$ , which is  $O\left(\frac{\log(1/\delta) + \log(1/\tau)}{\tau^2}\right)$ . Combining both cases, the expected running time of pSC is bounded by

$$O\left(\frac{\log(1/\delta) + \log(1/\max(\tau, \tau_{c}))}{\max(\tau, \tau_{c})^{2}}\right)$$
 (5)

# Proof of Theorem 2

**Theorem 12.** Let  $n = \frac{S}{\theta}$ , and  $\delta = O(\frac{1}{\log n})$ , then the running time of RdBin is bounded by  $O\left(\frac{\log(\log n/\tau)}{\tau^2} \times \log n\right)$ 

*Proof.* Let  $\delta = \frac{c}{\log n}$ , where c is considered to be a constant factor and c > 0. By Equation [5] we can obtain that for each random walk phase of RdBin, the expected running time is bounded by  $O\left(\frac{\log(\log n/\tau)}{\tau^2}\right)$ . As the random walk will terminate within  $\log n$  steps, and thus the overall running is bounded by  $O\left(\frac{\log(\log n/\tau)}{\tau^2} \times \log n\right)$ .

**Theorem 13.** Let  $p_v = q(f_v), \ \eta = \max(q^{-1}(p_v \pm \tau)),$ and  $\delta = O(\frac{1}{\log n})$ , then (i) the random walk of RdBin exits at the correct node with constant failure probability; and (ii) the inference error is bounded by  $\max(\frac{\theta}{S}, \eta)$ .

П

*Proof.* At each step of the random walk, the probability of moving in the wrong direction is  $O(\frac{1}{\log n})$ . As such, by union bound, the overall failure probability is O(1). If the random walk exits correctly, it must have either reached a leaf node or an interval  $[f_a, f_b]$  where  $|q(\hat{f} = \frac{f_a + f_b}{2}) - q(f_v)| < \tau$ , indicating that the minimum comparison threshold has been met. For the first case, since  $|f_a - f_b| = \frac{\theta}{S}$ , thus the inference error must be bounded by  $\frac{\theta}{S}$ . For the second case, it is clear that the error must be smaller than  $\max(|q^{-1}(p_v + \tau) - f_v|, |q^{-1}(p_v - \tau) - f_v|)$ .

# **B.3** Proof of Theorem 3

**Proof.** (**Theorem** 3) Consider party  $P_i$  and let  $S^i = \{s_1^i, s_2^i, ..., s_t^i\}$  to be the stake history of  $P_i$ . Given  $S_t = \{s_1, s_2, ..., s_t\}$  and  $S_t' = \{s_1', s_2', ..., s_t'\}$  to be a pair of  $(\alpha, u)$ -neighbor stake profiles. Let  $x, y \in \mathbb{R}$  such that  $x \le y \le x + \alpha$ , and by definition of  $(\alpha, u)$ -neighbor we have  $s_u = x$ ,  $s_u' = y$ . We also consider a given output  $o \subset Range(Lkg)$  and next. we compute the following probabilities (Note that we simplify the term  $Pr[Lkg^{(P_i,S)} = o]$  as a conditional probability  $Pr[o \mid S]$ )

$$\Pr\left[o \mid s_u^i = s_u\right] \le \max_{S_t} \Pr\left[o \mid S_t\right] \times \sum_{S_t \in S} \Pr\left[S_t \mid s_u^i = x\right]$$

$$\Pr\left[o\mid s_u^i = s_u'\right] \geq \min_{S_t'} \Pr\left[o\mid S_t'\right] \times \sum_{S_t' \subset \mathbb{S}} \Pr\left[S_t'\mid s_u^i = y\right]$$

Note that any stake dataset  $S \subset \mathbb{S}$  is a  $(\alpha, u)$ -neighbor of itself. Thus it holds that  $\sum_{S_t \subset \mathbb{S}} \Pr\left[S_t \mid s_u^i = x\right] \leq \sum_{S_t' \subset \mathbb{S}} \Pr\left[S_t' \mid s_u^i = y\right]$ , since for any  $S_t$  there exists at least one  $(\alpha, u)$ -neighbor  $S_t'$ . In addition, by Definition 5.3 we know that  $\max_{S_t} \Pr\left[o \mid S_t'\right] \leq e^{\omega(t-u)} \times \min_{S_t'} \Pr\left[o \mid S_t'\right]$ . Thus by combing the two facts, we can conclude that:

$$\frac{\Pr\left[o \mid s_{u}^{i} = s_{u}\right]}{\Pr\left[o \mid s_{u}^{i} = s_{u}^{i}\right]} \leq \frac{\sum_{S_{t} \subset \mathbb{S}} \Pr\left[S_{t} \mid s_{u}^{i} = x\right]}{\sum_{S_{t}^{\prime} \subset \mathbb{S}} \Pr\left[S_{t}^{\prime} \mid s_{u}^{i} = y\right]} \times \frac{\max_{S_{t}} \Pr\left[o \mid S_{t}\right]}{\min_{S_{t}^{\prime}} \Pr\left[o \mid S_{t}^{\prime}\right]}$$

By Bayes' theorem

$$\frac{\Pr\left[s_{u}^{i} = x \mid o\right]}{\Pr\left[s_{u}^{i} = y \mid o\right]} / \frac{\Pr\left[s_{u}^{i} = x\right]}{\Pr\left[s_{u}^{i} = y\right]}$$

$$= \frac{\Pr\left[s_{u}^{i} = x \mid o\right] \Pr\left[o\right]}{\Pr\left[s_{u}^{i} = x\right]} / \frac{\Pr\left[s_{u}^{i} = y \mid o\right] \Pr\left[o\right]}{\Pr\left[s_{u}^{i} = y\right]}$$

$$= \frac{\Pr\left[o \mid s_{u}^{i} = s_{u}\right]}{\Pr\left[o \mid s_{u}^{i} = s_{u}'\right]} \le e^{\omega(t-u)\times\epsilon} \left(\mathbf{Definition} \mathbf{5.1}\right)$$
(6)

## B.4 Proof of Theorem 4

*Proof.* (**Theorem** 4) We first recall the important notations: (i) REAL $_{\mathcal{A},\mathcal{E}}^{\pi}$ , denotes the transcript of real world

protocol  $(\pi)$  execution between the honest party and the adversary  $\mathcal R$  involving environment  $\mathcal E$ ; (ii) IDEAL $_{S,\mathcal E}^{\mathcal F}$ , denotes the transcript generated by simulator a  $\mathcal S$  with input of a leakage profile Lkg who corrupts the same parties as  $\mathcal R$  and interacts only with a trusted functionality  $\mathcal F$  and the environment  $\mathcal E$ . We consider a set of parties  $P=\{P_1,P_2,...,P_n\}$ , and the environment  $\mathcal E$  can activate stake distribution,  $\mathbb S=\{S^i\}_{1\leq i\leq n}$ , for all parties across every time slot. (In fact, the environment  $\mathcal E$  can activate every transaction for each party, which essentially governs stake distribution among parties over time.) By definition, as  $\prod$  UC-emulates ledger maintenance command, thus  $\forall$   $\mathcal E$ , there exists a simulator  $\mathcal S_{\text{LM}}$  such that:

$$\mathsf{IDEAL}^{\mathcal{F}_{\mathsf{LM}}}_{\mathcal{S}_{\mathsf{LM}}(\phi(\mathbb{S})),\mathcal{E}} \approx_{\mathsf{ind}} \mathsf{REAL}^{\Pi(\mathbb{S})}_{\mathcal{A},\mathcal{E}} \tag{7}$$

Next, we prove that with the same simulator  $\mathcal{S}_{LM}$ , one can also simulate indistinguishable transcripts against protocol  $\prod_{LM^*}$ . We prove this by contradiction. Assume there exists an environment  $\mathcal{E}$  and a distorted stake distribution  $\tilde{\mathbb{S}}_a \leftarrow f(\mathbb{S}_a)$  such that

$$\mathsf{IDEAL}^{\mathcal{T}_\mathsf{LM}}_{\mathcal{S}_\mathsf{LM}(\phi(\tilde{\mathbb{S}}_a))}, \mathcal{E} \not\approx_\mathsf{ind} \mathsf{REAL}^{\prod_{\mathsf{LM}^*}(\mathbb{S}_a)}_{\mathcal{A}, \mathcal{E}} \tag{8}$$

where  $\prod_{\mathsf{LM}^*}(\mathbb{S}_a)$  denotes the execution of protocol  $\prod_{\mathsf{LM}^*}$  under the stake distribution  $\mathbb{S}_a$ . As  $\prod_{\mathsf{LM}^*}$  only modifies the inputs rather than the interface with the environment or the components of  $\prod$ , the equivalence  $\mathsf{REAL}_{\mathcal{A},\mathcal{E}}^{\prod_{\mathsf{LM}^*}(\mathbb{S}_a)} \approx_{\mathsf{ind}} \mathsf{REAL}_{\mathcal{A},\mathcal{E}}^{\prod_{(\tilde{\mathbb{S}}_a)}}$  becomes evident. In other words, we keep the environment unchanged but only change the stake distribution from  $\mathbb{S}_a$  to  $\tilde{\mathbb{S}}_a$ , then the execution transcripts of  $\prod_{\mathsf{LM}^*}(\mathbb{S})$  and  $\prod_{(\tilde{\mathbb{S}})}$  should be indistinguishable. And by the assumption Eq  $\boxed{\mathbb{S}}$ , we can trivially obtain that

$$\mathsf{IDEAL}^{\mathcal{T}_\mathsf{LM}}_{\mathcal{S}_\mathsf{LM}(\phi(\tilde{\mathbb{S}}_a))}, \mathcal{E} \not\approx_{\mathsf{ind}} \mathsf{REAL}^{\prod(\tilde{\mathbb{S}}_a)}_{\mathcal{A}, \mathcal{E}}$$

This implies that if  $\mathcal{E}$  activates the stake distribution as  $\tilde{\mathbb{S}}_a$ , then simulator  $\mathcal{S}_{LM}$  is unable to generate indistinguishable transcripts against the real protocol executions of  $\Pi$ , which evidently contradicts Equation  $\Pi$ . Consequently, we can deduce that no such environment exists, and the assumption (Eq  $\Pi$ ) is incorrect. Therefore, for any stake distribution activation, the same simulator  $\mathcal{S}_{LM}$  can also produce indistinguishable transcripts against the real execution of  $\Pi_{LM^*}$  with leakage  $\phi(\tilde{\mathbb{S}})$ .

## B.5 Proof of Theorem 5

*Proof.* (**Theorem** [5]) Let  $S_t = \{s_j\}_{1 \le j \le t}$  and  $S'_t = \{s'_j\}_{1 \le j \le t}$  be any  $(\alpha, u)$ -neighboring stake assignments for an honest party P. We abstract the distortion at each time j as  $\mathcal{M}_j(s)$ , which takes an input stake, s, and outputs  $s + \text{Lap}(\frac{\alpha}{\epsilon})$  if  $j \mod T = 0$ , and  $\mathcal{M}_{j-1}(s)$ 

otherwise. For ease of notation, we write the conditional probability  $\Pr \left[ \mathcal{M}_i(s_i) = o_i \mid \forall_{1 \le k \le j-1} \mathcal{M}_k(s_k) = o_k \right]$ as  $\Pr [\mathcal{M}_i(s_i) = o_i \mid *]$ , and compute the following term

$$\frac{\Pr\left[\tilde{S}_{t}=\mathbf{o}\right]}{\Pr\left[\tilde{S}_{t}'=\mathbf{o}\right]} \leq \prod_{j=1}^{u-1} \mathbf{1} \times \prod_{j=u \land j=kT, k \in \mathbb{N}^{+}}^{t} e^{\epsilon} \leq e^{\left\lfloor \frac{t-u}{T} \right\rfloor \times \epsilon}$$

In addition, knowing that  $Lkg = f(\tilde{S}_t)$  is a probabilistic function related to the noisy stake assignment. Thus, by post-processing theorem of DP [26],  $\Pr[\mathsf{Lkg} = \mathbf{o}] / \Pr[\mathsf{Lkg}' = \mathbf{o}] \le e^{\lfloor \frac{t-u}{T} \rfloor \times \epsilon}$ . In general, the total privacy loss for stake values at any time u is subject to k-fold composition theorem of DP mechanisms [26], where k denotes the total number of noisy stake releases from time u up to the current moment, which is  $\frac{t-u}{T}$ .

#### **B.6 Proof of Theorem 6**

*Proof.* (**Theorem** 6) Let  $S_t = \{s_j\}_{1 \le j \le t}$  to be the stake profiles of an honest party P, and  $tx_t = \{tx_j\}_{1 \le j \le t}$  to be the corresponding transaction outcomes at each time, i.e.,  $tx_i = s_i - s_{i-1}$ . Let  $M(tx_t)$  be a mechanism that processes the noisy p-sums based on  $tx_t$ , and  $\tilde{S}_t \leftarrow f(M(tx_t))$  is an algorithm that aggregates the output of M to derive the noisy stake at each time. We say that, for any pair of  $(\alpha, u)$ -neighboring stake profiles  $S_t$ , and  $S'_t$ , with corresponding  $tx_t$ ,  $tx'_t$ . It holds that if for any  $\mathbf{o} \subseteq Range(M)$ ,

$$\frac{\Pr\left[M(\mathsf{tx}_t) = \mathbf{0}\right]}{\Pr\left[M'(\mathsf{tx}_t) = \mathbf{0}\right]} \le e^{\omega(t-u)\times\epsilon} \tag{9}$$

hence the released noisy stake  $\tilde{S}_t \leftarrow f(M(\mathsf{tx}_t))$  satisfies Definition 5.3 under the same decay function  $\omega$ , given f only conduct post-processing operations. By the theory of the Laplace mechanism,  $\frac{\Pr[M(\mathsf{t}_t)=0]}{\Pr[M'(\mathsf{t}_t)=0]} \le e^{\epsilon}$ , when t=u. Moreover, when  $0 < t-u \le L$ , then the privacy loss is subject to phase two distortion (privacy loss due to noisy interval tree generation). Note that  $tx_u$  can be used to generate at most  $\log_2(\frac{(t-u)}{T})$  many noisy nodes (p-sums), thus the phase two privacy loss is bounded by

$$\Pr[\mathsf{Lkg} = \mathbf{o}] / \Pr[\mathsf{Lkg}' = \mathbf{o}] \le e^{\log_2(\frac{t-u}{T}) \times \epsilon}$$

In addition, when t - u > L, then the mechanism incurs privacy loss for both phase one and phase two distortion. By Theorem 5 the total phase one privacy loss is bounded by  $\lfloor \frac{t-u}{L} \rfloor \times \epsilon$ . For phase two, the max privacy loss is bounded by  $e^{\log_2(\frac{L}{T})\times\epsilon}$ . As such, we may conclude

$$\Pr\left[\mathsf{Lkg} = \mathbf{o}\right] / \Pr\left[\mathsf{Lkg'} = \mathbf{o}\right] \le e^{\left(\left\lfloor \frac{t-u}{L} \right\rfloor + \log_2\left(\frac{t-u}{T}\right)\right) \times \epsilon}$$

As summary, the privacy decay function to be

$$\omega(t-u) \le \begin{cases} \lfloor \frac{t-u}{L} \rfloor + \log_2(\frac{L}{T}), & \text{if } (t-u) > L\\ \log_2(\frac{t-u}{T}), & \text{otherwise} \end{cases}$$
(10)

## **Proof of Theorem 7**

**Lemma 14.** Given n independent and identically distributed (i.i.d.) Laplace random variables  $X_1, X_2, ..., X_n$ drawn from Lap $(\frac{\Delta}{\epsilon})$ . Let  $X = \sum_{i=1}^{n} X_i$ ,  $0 < \alpha \le n \frac{\Delta}{\epsilon}$ , then:

$$\Pr[X \ge \alpha] \le e^{\left(\frac{-\alpha^2 \Delta^2}{4n\epsilon^2}\right)}$$

*Proof.* Please refer to the proof to Lemma 12.2 in [26] or proof to Theorem 6 in [55]

*Proof.* (**Theorem** 7) When adopting stake distortion, the inference errors consist of two parts: (i) the error due to injected DP noises and (ii) the error caused by the inference algorithm. By Theorem 2, the part (ii) noise is bounded by  $\eta$ . By Lemma 14, and let  $e^{(\frac{-\alpha^2 \Delta^2}{4n\epsilon^2})} = \beta$ , and take log on both sides, one can obtain that with probability at least  $1 - \beta$ , for any  $\beta > 0$ . X is bounded by  $2\frac{\Delta}{\epsilon} \sqrt{n \log \frac{1}{\beta}}$ . Hence, we can conclude that when  $\mathcal{F}_{\mathsf{Timer}}$  and  $\mathcal{F}_{\mathsf{Bin}}$  are applied, the part (i) errors are bounded by  $O(\frac{\alpha}{\epsilon})$ , and  $O(\frac{\sqrt{L\alpha}}{\epsilon})$ , respectively.

#### **Proof of Theorem 8 B.8**

*Proof.* (**Theorem 8**) Let  $X_a \le \frac{1-\gamma}{3}S$ ,  $X_h \ge \frac{2+\gamma}{3}S$  to be stake controlled by the adversary and the honest parties, respectively. We denote the distorted stake as  $\tilde{X}_a \leftarrow X_a + Y^{n_a}$ , and  $\tilde{X}_h \leftarrow X_h + Y^{n_h}$ , respectively for adversary and honest, where  $Y^x$  is the summation of x *i.i.d.* Laplace random variables, and  $n_a$ ,  $n_b$  denotes the number of honest and malicious parties, respectively, such that  $n_a + n_b = n$ . To ensure that the corrupted stake after distortion is bounded by  $\frac{1}{3}$ , we need to have  $\tilde{X}_h - 2\tilde{X}_a > 0$ with high probability. Knowing that:

$$\tilde{X}_h - 2\tilde{X}_a = (X_h - 2X_a) + Y^{n_h} - 2Y^{n_a} = \gamma S + Y^{n_h} - 2Y^{n_a}$$
(11)

Since the random variables are symmetric, without loss of generality, we may assume the following

$$\Pr\left[\tilde{X}_h - 2\tilde{X}_a < 0\right] = \Pr\left[Y^{n_h} + 2Y^{n_a} > \gamma S\right] \tag{12}$$

By Lemma 12.2 in [26], we can obtain that for  $\beta \in (0,1)$ 

$$\Pr\left[Y^{n_h} + 2Y^{n_a} > \sqrt{\log \frac{2}{\beta}} \frac{\alpha \left(\sqrt{n_h} + 2\sqrt{n_a}\right)}{\epsilon S}\right] \le \beta \quad (13)$$

As  $n_h + n_a < \frac{S}{n}$  and  $n_h > 2n_a$ , then  $\sqrt{n_h} + 2\sqrt{n_a} < \frac{S}{n_h}$  $1.74\sqrt{\frac{S}{v}}$ . Thus, by setting  $\gamma \ge \frac{1.74\alpha}{\epsilon S}\sqrt{\frac{S}{v}\log\frac{1}{\beta}}$ , one can conclude that with probability at least  $1 - \beta$ , the adversary controlled stake after distortion does not exceed  $\frac{1}{3}$ .